

**Анотація.** У статті описується розроблений алгоритм пошуку залежностей у підмножинах об'єктів досліджуваного набору даних. Даний алгоритм дозволяє ефективно виявляти асоціативні залежності за заданими критеріями якості. Застосування розробленого методу виявлення залежностей в даних можливе в багатьох предметних галузях та дозволяє виявити нові закономірності в даних, що покращує роботу фахівців та якість прийнятих ними рішень.

**Ключові слова:** залежність даних, набір даних, алгоритм, аналіз даних.

**Аннотация.** В статье описан разработанный алгоритм поиска зависимостей в подмножествах объектов исследуемого набора данных. Данный алгоритм позволяет эффективно выявлять ассоциативные зависимости в данных по заданным критериям качества. Использование разработанного метода поиска зависимостей в данных возможно в многих предметных областях и позволяет выявить новые закономерности в данных, что улучшает работу специалистов и качество принятых ими решений.

**Ключевые слова:** зависимость данных, набор данных, алгоритм, анализ данных.

**Abstract.** The worked out algorithm of data search dependencies in the object subsets of investigated dataset is described in the article. This algorithm allows effectively reveal associative dependencies based on specified quality criteria. Application of developed data dependency detection method is possible in many specializations and allows finding new data patterns which improves work and decisions quality of specialists.

**Keywords:** data dependency, dataset, algorithm, data analysis.

## 1. Вступ

Аналіз результатів будь-яких досліджень чи спостережень є невід'ємною частиною процесу дослідження предметної області. Напряму аналізу інформації швидко зростає і поширюється на нові й нові галузі науки і техніки. Водночас з галузевим поширенням методи аналізу даних застосовуються до все більших наборів даних. Це потребує розробки нових підходів та алгоритмів у даній галузі. Хоча на сьогодні наявні потужні методи та засоби аналізу даних у деяких вузьких галузях, наприклад, CLASSIFI (Department of Pathology, UT Southwestern Medical Center) [1], BiNGO (Department of Plant Systems Biology, VIB/Ghent University) [2] та EASE (National Institute of Allergy and Infectious Diseases) [3], проте більшість зібраної інформації не аналізується належним чином перш за все через недостатні обчислювальні потужності для наявних алгоритмів і відсутність ефективніших методів аналізу цих даних. Внаслідок цього багато залежностей предметних областей залишаються невідомими експертам та аналітикам і це не дозволяє приймати правильні рішення у проблемних ситуаціях, що приносить великі збитки як економічні, так і екологічні, медичні та ін. Тому розробка алгоритму, застосовного до широкого спектра схем даних, є надзвичайно актуальною задачею.

Метою даного дослідження є розробка ефективного алгоритму аналізу широкого спектра даних, що дозволить фахівцям предметних областей краще зрозуміти залежності в цих галузях та приймати правильні рішення у проблемних ситуаціях.

## 2. Основна частина

Аналіз великих обсягів даних потребує виявлення груп атрибутів, що утворюють функціональні залежності. Проте в реальних умовах значно частіше зустрічаються набори даних, в

яких важливі залежності визначені тільки на деякій підмножині значень групи ключових атрибутів. Називатимемо такі залежності частковими функціональними залежностями (ЧФЗ).

Тобто, часткова функціональна залежність – це ФЗ в деякій селекції основного відношення.

$$F_p : K = \{a_i\}, a_i \in A, D = \{a_j\}, a_j \in A, R' \subset R : K \rightarrow D | R'. \quad (1)$$

Багато залежностей також мають не чітко детермінований характер. Називатимемо їх імовірнісними продукційними залежностями (ІПЗ).

Імовірнісна продукційна залежність – це продукційне правило в селекції основного відношення, яке справджується для значущої кількості об'єктів цієї селекції. Поріг значущості повинен визначатись експертним шляхом або виходячи з розрахунків імовірності помилкового виділення цієї залежності.

$$F_I : K = \{a_i\}, a_i \in A, D = \{a_j\}, a_j \in A, : P(k \in K \rightarrow d \in D) = p. \quad (2)$$

Тут  $k$  та  $d$  – кортежі значень деяких груп атрибутів  $K$  та  $D$  відповідно.

Основним показником достовірності такої залежності є відношення кількості об'єктів, для яких має місце така ФЗ, до кількості об'єктів у селекції:

$$P(F_I) = \frac{|\sigma_{k \in K \wedge d \in D}(R)|}{|\sigma_{k \in K}(R)|}. \quad (3)$$

Позначатимемо ІЗ  $F_I$  з використанням цієї величини як імовірність виконання продукційного правила:

$$P(k \in K \rightarrow d \in D) = p. \quad (4)$$

При потребі неважко розширити поняття до імовірності виконання довільних продукційних правил. Але пошук складних продукційних правил у великих наборах даних є дуже складною і обчислювально важкою операцією, тому вони використовуються рідко.

Введення поняття імовірнісних залежностей, на перший погляд, надлишкове, оскільки з множини об'єктів, на яких визначена така залежність, завжди можна вибрати підмножину, на якій визначена чітка часткова функціональна залежність. Але за умови, що значення групи ключових атрибутів такої залежності мають змістовне сортування у предметній області, опис залежностей може бути здійснений значно простіше та інформативніше. Для наочності розглянемо приклад:

Приклад 1.

Таблиця 1. Залежність типу вагона, в якому їде людина, від дальності поїздки

і – номер опитаного	L – дальність поїздки, x100км	T – клас вагона
1	1	3
2	1	2
3	1	3
4	3	1
5	2	1
6	1	3
7	2	1
8	2	2
9	3	1

З табл. 1 можна зробити висновок.

Такі ЧФЗ:

$$i \rightarrow (T, L) \{1..9\},$$

$$L \rightarrow T \{3\},$$

$$T \rightarrow L \{3\},$$

$$(L, T) \rightarrow i \{(1;2), (2;2)\}.$$

Розглядаються тільки максимально розширені ЧФЗ, тобто, наприклад, залежність  $i \rightarrow (T, L) \{1,2,3\}$  не береться

до уваги, оскільки перша ЧФ з вищенаведених повністю визначає її.

Ці залежності справджуються для всіх об'єктів, які мають вказану комбінацію визначальних атрибутів залежності, але інформативність таких залежностей часто мала через надмірний детермінізм – будь-який виняток зруйнує залежність і будь-який випадковий шум даних може утворити нову залежність.

Розглянемо деякі ПЗ з попереднього прикладу:

$$P(i = 1 \rightarrow (T; L) = (1; 3)) = 1,$$

$$P(L = 1 \rightarrow T = 3) = 0,75,$$

$$P(L \in \{2, 3\} \rightarrow T \in \{1, 2\}) = 1,$$

$$P(L \in \{2, 3\} \rightarrow T = 1) = 0,8,$$

$$P(T = 3 \rightarrow L = 1) = 1,$$

$$P(T = 1 \rightarrow L = 3) = 0,5.$$

Такі дані значно інформативніші і корисніші для аналізу. Вони захищені від шуму та випадкових помилок і навіть дозволяють виявити ці помилки. Тому пошук саме таких залежностей і є завданням сучасного аналізу даних. Задачу пошуку мінімального покриття часто використовують в аналізі даних [4].

Ще одним важливим параметром ПЗ є її достовірність – імовірність відхилення величини  $P(F_i)$  не більше, ніж на задане значення  $\delta$ :

$$P(|P(F_i) - \delta| \leq P^*(F_i)), \quad (5)$$

де  $P^*(F_i)$  – справжнє значення виконання залежності  $F_i$ .

Обчислення  $P^*(F_i)$  напряму в більшості випадків неможливе, оскільки відсутня інформація про всі можливі комбінації значень атрибутів об'єктів відношення. Але можна оцінити достовірність залежності, виходячи з кількості об'єктів, для яких залежність виконується, комбінацій задіяних атрибутів та інших параметрів.

Легко зауважити, що ЧФ є лише частковим випадком ПЗ при  $P(F_i) = 1$ . Тому завдання пошуку залежностей у даних можна звести до пошуку ПЗ.

У даній роботі пропонується метод виявлення ПЗ в наборах даних, опираючись на задані параметри імовірності виконання та достовірності цих залежностей.

Розглянемо для початку пошук простих залежностей типу

$$P(s = s_1 \rightarrow t = t_1) \geq p_0, \quad (6)$$

де  $p_0$  – порогове значення імовірності виконання шуканих залежностей,  $s, t$  – деякі атрибути відношення, в якому здійснюється пошук залежностей.

Для обчислення необхідних імовірностей та пошуку таких ПЗ необхідно перебрати всі пари атрибутів  $(s; t)$ ,  $s \in A, t \in A$  і для кожної комбінації здійснити прохід по об'єктах відношення, заповнюючи для кожного значення атрибута  $s$  кількість його повторень в об'єктах та словник відповідних значень атрибута  $t$  із вказанням кількості повторів. Після таких операцій усі залежності виду (6) при  $p_0 = 0$  будуть збережені в описаній структурі даних. Відбір потрібних ПЗ з цієї структури для довільного  $p_0$  є тривіальною задачею фільтрації. Складність описаного алгоритму –  $O(m^2 \cdot n)$ , де  $m$  – кількість атрибутів відношення,  $n$  – кількість об'єктів у відношенні.

Хоча складність вищеописаного алгоритму і залежить не лінійно від кількості атрибутів відношення, проте переважно навіть дуже складні дані мають не більше декількох сотень атрибутів. Окрім того, алгоритм має добрі дані розпаралелювання, тобто, його можна елементарно модифікувати для виконання на багатопроцесорних машинах, кластерах чи ґрідах.

Розширимо задачу (6) до пошуку залежностей виду

$$P(s = s_1 \rightarrow t \in T) \geq p_0. \quad (7)$$

Маючи дані про залежності виду (6), можна розв'язати дану задачу, аналізуючи лише підгрупу залежностей виду (6), що має спільну частину  $s = s_1$ . Це легко довести, виходячи з визначення ІЗ та формули (3): будь-який об'єкт, для якого діє залежність  $P(s = s_2 \rightarrow q \in Q) \geq p$ , не виконує залежність (7) і, відповідно, не змінює значення  $P(s = s_1 \rightarrow t \in T)$ . З іншої сторони, додавання чи вилучення об'єкта зі значенням  $s = s_1$  неодмінно впливає на імовірність  $P(s = s_1 \rightarrow t \in T)$ , якщо це значення більше нуля. Отже, при пошуку залежностей виду (7) необхідно і достатньо розглянути всі залежності виду (6), у яких умовна частина предиката співпадає з шуканою залежністю.

Щодо імовірності виконання залежності (7) на заданому наборі даних, то з (3) випливає, що

$$\begin{aligned} P(s = s_1 \rightarrow t \in \{t_1, t_2\}) &= \frac{|\sigma_{t \in \{t_1, t_2\}}(R)|}{|\sigma_{s=s_1}(R)|} = \frac{|\sigma_{t=t_1}(R)| + |\sigma_{t=t_2}(R)|}{|\sigma_{s=s_1}(R)|} = \\ &= P(s = s_1 \rightarrow t = t_1) + P(s = s_1 \rightarrow t = t_2). \end{aligned} \quad (8)$$

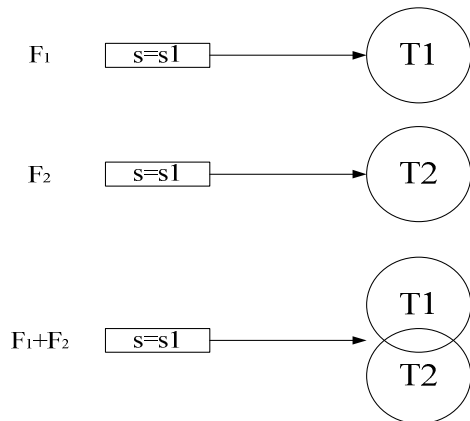


Рис. 1. Графічне представлення об'єднання ІЗ типу (7)

Тобто, об'єднання залежностей виду (6) у залежності від виду (7) збільшить імовірність виконання нової залежності. Для більш загального випадку об'єднання двох залежностей виду (7) формула обчислення імовірності виконання агрегативної залежності буде трохи складнішою, оскільки множини значень результуючої частини предиката можуть перекриватись (рис. 1).

Відповідно, виразимо з (3) імовірність виконання такої ІЗ:

$$\begin{aligned} P(s = s_1 \rightarrow t \in T_1 \cup T_2) &= \frac{|\sigma_{t \in T_1 \cup T_2}(R)|}{|\sigma_{s=s_1}(R)|} = \\ &= \frac{|\sigma_{t \in T_1}(R)| + |\sigma_{t \in T_2}(R)| - |\sigma_{t \in T_1 \cap T_2}(R)|}{|\sigma_{s=s_1}(R)|}. \end{aligned}$$

Використовуючи зворотне перетворення формули (3), легко помітити, що

$$\begin{aligned} \frac{|\sigma_{t \in T_1}(R)|}{|\sigma_{s=s_1}(R)|} &= P(s = s_1 \rightarrow t \in T_1), \\ \frac{|\sigma_{t \in T_2}(R)|}{|\sigma_{s=s_1}(R)|} &= P(s = s_1 \rightarrow t \in T_2), \end{aligned}$$

$$\frac{|\sigma_{t \in T_1 \cap T_2}(R)|}{|\sigma_{s=s_1}(R)|} = P(s = s_1 \rightarrow t \in T_1 \cap T_2),$$

і утворене відношення можна переписати простіше:

$$P(s = s_1 \rightarrow t \in T_1 \cup T_2) = P(s = s_1 \rightarrow t \in T_1) + P(s = s_1 \rightarrow t \in T_2) - P(s = s_1 \rightarrow t \in T_1 \cap T_2). \quad (9)$$

Враховуючи обчислені значення імовірностей всіх атомарних залежностей (6), кон'юнкція та диз'юнкція довільних множин значень атрибута Т не потребує нових складних обчислень для знаходження імовірності виконання об'єднаної залежності.

Ще одним варіантом об'єднання ПЗ виду (6) є залежності

$$P(s \in S \rightarrow t = t_1) \geq p_0. \quad (10)$$

Обчислення імовірності такої агрегативної залежності трохи складніше, ніж у попередньому випадку і неможливе, якщо спиратись тільки на значення  $P(s = s_i \rightarrow t = t_1)$ . Для цього потрібно залишати імовірності залежностей виду (6) у вигляді звичайних дробів – відношення кількості об'єктів, для яких залежність виконується до загальної кількості об'єктів, відібраних умовною частиною предиката ПЗ, як і було запропоновано в алгоритмі пошуку ПЗ виду (6). Тоді, виходячи з формули (3),

$$P(s \in \{s_i\} \rightarrow t = t_1) = \frac{\sum_i |s = s_i \wedge t = t_1|}{\sum_i |s = s_i|}. \quad (11)$$

І знову ж, використовуючи вищеописану структуру даних, ця формула обчислюється без додаткових проходів по таблиці.

Розглянемо тепер загальний вигляд шуканих залежностей:

$$P(s \in S \rightarrow t \in T) \geq p_0. \quad (12)$$

Обчислення імовірності виконання такої залежності ґрунтується на попередніх міркуваннях та можливості розкладу такої залежності на складові ПЗ типу (11):

$$P(s \in S \rightarrow t \in T) = \sum_{t_i \in T} P(s \in S \rightarrow t = t_i) = \sum_{t_i \in T} \frac{\sum_j |s = s_j \wedge t = t_i|}{\sum_j |s = s_j|}. \quad (13)$$

Позначення, що використовувались до цих пір у статті, досить громіздкі, тому введемо алгебру залежностей, визначену над четвірками

$$(S; T; |s \in S \wedge t \in T|; |s \in S|) = F_t, \quad (14)$$

операціями об'єднання (+) та розщеплення (-) залежностей. Для додаткового спрощення позначень введемо операцію проєкції трійки  $F_t$  на її атрибути (назвемо їх відповідно Pr (Predicate), PrS, PrT (дочірні атрибути предиката  $s \in S \rightarrow t \in T$  – множини S та T), PN(Probability nominator), PD (Probability Denominator) і позначатимемо операцію квадратними дужками, наприклад,  $F_t[\text{Pr}]$ ,  $F_t[\text{PrS}]$ .

Для типів ПЗ, описаних формулами (6), (7), (10), (12), введемо позначення відповідно  $F^A$ ,  $F^B$ ,  $F^C$ ,  $F^D$ . Тоді операції об'єднання та розщеплення над залежностями цих типів будуть визначені таким чином:

$$F_1^A + F_2^A = \begin{cases} F_1^A, & \text{якщо } F_1^A[\text{Pr } S] = F_2^A[\text{Pr } S] \wedge F_1^A[\text{Pr } T] = F_2^A[\text{Pr } T], \\ F^B = (F_1^A[\text{Pr } S]; F_1^A[\text{Pr } T] \cup F_2^A[\text{Pr } T]; F_1^A[\text{Pr } T] + F_2^A[\text{Pr } T]; F_1^A[\text{Pr } T]), \\ & \text{якщо } F_1^A[\text{Pr } S] = F_2^A[\text{Pr } S] \wedge F_1^A[\text{Pr } T] \neq F_2^A[\text{Pr } T], \\ F^C = (F_1^A[\text{Pr } S] \cup F_2^A[\text{Pr } S]; F_1^A[\text{Pr } T]; F_1^A[\text{Pr } T] + F_2^A[\text{Pr } T]; F_1^A[\text{Pr } T] + F_2^A[\text{Pr } T]), \\ & \text{якщо } F_1^A[\text{Pr } S] \neq F_2^A[\text{Pr } S] \wedge F_1^A[\text{Pr } T] = F_2^A[\text{Pr } T], \\ F^D = (F_1^A[\text{Pr } S] \cup F_2^A[\text{Pr } S]; F_1^A[\text{Pr } T] \cup F_2^A[\text{Pr } T]; F_1^A[\text{Pr } T] + F_2^A[\text{Pr } T]; F_1^A[\text{Pr } T] + F_2^A[\text{Pr } T]), \\ & \text{якщо } (F_1^A[\text{Pr } S] \neq F_2^A[\text{Pr } S]) \wedge (F_1^A[\text{Pr } T] \neq F_2^A[\text{Pr } T]), \end{cases} \quad (15)$$

$$F^B + F^A = \begin{cases} F^B, & \text{якщо } F^B[\text{Pr } S] = F^A[\text{Pr } S] \wedge F^A[\text{Pr } T] \in F^B[\text{Pr } T], \\ F_1^B = (F^B[\text{Pr } S]; F^B[\text{Pr } T] \cup F^A[\text{Pr } T]; F^B[\text{Pr } T] + F^A[\text{Pr } T]; F^B[\text{Pr } T]), \\ & \text{якщо } F^B[\text{Pr } S] = F^A[\text{Pr } S] \wedge F^A[\text{Pr } T] \notin F^B[\text{Pr } T], \\ F^D = (F^B[\text{Pr } S] \cup F^A[\text{Pr } S]; F^B[\text{Pr } T] \cup F^A[\text{Pr } T]; F^B[\text{Pr } T] + F^A[\text{Pr } T]; F^B[\text{Pr } T] + F^A[\text{Pr } T]), \\ & \text{якщо } F^B[\text{Pr } S] \neq F^A[\text{Pr } S], \end{cases} \quad (16)$$

$$F^C + F^A = \begin{cases} F^C, & \text{якщо } F^A[\text{Pr } S] \in F^C[\text{Pr } S] \wedge F^C[\text{Pr } T] = F^A[\text{Pr } T], \\ F_1^C = (F^C[\text{Pr } S] \cup F^A[\text{Pr } S]; F^C[\text{Pr } T]; F^C[\text{Pr } T] + F^A[\text{Pr } T]; F^C[\text{Pr } T] + F^A[\text{Pr } T]), \\ & \text{якщо } F^A[\text{Pr } S] \notin F^C[\text{Pr } S] \wedge F^C[\text{Pr } T] = F^A[\text{Pr } T], \\ F_1^D = (F^C[\text{Pr } S]; F^C[\text{Pr } T] \cup F^A[\text{Pr } T]; F^C[\text{Pr } T] + F^A[\text{Pr } T]; F^C[\text{Pr } T]), \\ & \text{якщо } (F^A[\text{Pr } S] \in F^C[\text{Pr } S]) \wedge (F^C[\text{Pr } T] \neq F^A[\text{Pr } T]), \\ F_2^D = (F^C[\text{Pr } S] \cup F^A[\text{Pr } S]; F^C[\text{Pr } T] \cup F^A[\text{Pr } T]; F^C[\text{Pr } T] + F^A[\text{Pr } T]; F^C[\text{Pr } T] + F^A[\text{Pr } T]), \\ & \text{якщо } (F^A[\text{Pr } S] \notin F^C[\text{Pr } S]) \wedge (F^C[\text{Pr } T] \neq F^A[\text{Pr } T]), \end{cases} \quad (17)$$

$$F^D + F^A = \begin{cases} F^D, & \text{якщо } (F^A[\text{Pr } S] \in F^D[\text{Pr } S]) \wedge (F^A[\text{Pr } T] \in F^D[\text{Pr } T]), \\ F_1^D = (F^D[\text{Pr } S]; F^D[\text{Pr } T] \cup F^A[\text{Pr } T]; F^D[\text{Pr } T] + F^A[\text{Pr } T]; F^D[\text{Pr } T]), \\ & \text{якщо } (F^A[\text{Pr } S] \in F^D[\text{Pr } S]) \wedge (F^A[\text{Pr } T] \notin F^D[\text{Pr } T]), \\ F_2^D = (F^D[\text{Pr } S] \cup F^A[\text{Pr } S]; F^D[\text{Pr } T] \cup F^A[\text{Pr } T]; F^D[\text{Pr } T] + F^A[\text{Pr } T]; F^D[\text{Pr } T] + F^A[\text{Pr } T]), \\ & \text{якщо } (F^A[\text{Pr } S] \notin F^D[\text{Pr } S]). \end{cases} \quad (18)$$

Залежності типів  $F^B$ ,  $F^C$ ,  $F^D$  можна представити у вигляді суми ІІЗ типу  $F^A$ , виразивши їх параметри зворотними перетвореннями формул (15)–(18). Тому наводити формули обчислення сум  $F^B + F^C$ ,  $F^D + F^D$  та ін. тут немає сенсу.

Використовуючи (15)–(18), можна елементарно довести, що операція об'єднання залежностей має властивості комутативності та асоціативності на множині ІІЗ.

Використовуючи (3) та (15)–(18), визначимо операцію розщеплення ІІЗ як обернену до операції об'єднання:

$$F_1 + F_2 = F_3 \Rightarrow F_3 - F_2 = F_1. \quad (19)$$

Таким чином, операції розширення та звуження множин значень умови та результату продукційного правила ІІЗ дозволяють гнучко модифікувати ці залежності без потреби проведення додаткових складних обчислень. Необхідно тільки провести попередній аналіз простих залежностей виду (6) вищеописаним алгоритмом.

Наведені операції та алгоритм дозволяють виявити область визначення ПЗ та знайти залежності за заданою імовірністю виконання. Недоліком залишається знаходження залежностей тільки виду (12) (вони включають (6), (7) та (10)), тобто ПЗ, що мають лише один атрибут в умовній частині предиката та один – у результуючій. Тобто, неможливий пошук залежностей

$$P(s_1 \in S_1 \wedge s_2 \in S_2 \wedge \dots \wedge s_l \in S_l \rightarrow t_1 \in T_1 \wedge t_2 \in T_2 \wedge \dots \wedge t_k \in T_k) \geq p_0. \quad (20)$$

Для забезпечення пошуку таких ПЗ введемо дві нові операції: введення нового атрибуту в умовну частину предиката залежності та введення нового атрибуту в результуючу частину предиката. Позначимо їх  $\otimes$  та  $\oplus$  відповідно.

Для залежностей (20) розширимо поняття операцій  $F_l[\text{Pr}]$ ,  $F_l[\text{Pr } S]$ . Вони повертають множину значень кортежу атрибутів умови та результуючої частини предиката відповідно. Ця множина значень є декартовим добутком множин допустимих значень окремих атрибутів кортежу:

$$\begin{aligned} F_l[\text{Pr } S] &= S_1 \times S_2 \times \dots \times S_l, \\ F_l[\text{Pr } T] &= T_1 \times T_2 \times \dots \times T_k. \end{aligned} \quad (21)$$

Операції з ПЗ тепер визначимо як

$$F_l \oplus (x \in X) = F_l' = \left( F_l[\text{Pr } S]; F_l[\text{Pr } T] \times X; \left| \sigma_{x \in X} \left( \sigma_{k \in F_l[\text{Pr } S]}(R) \right) \right|; \left| \sigma_{k \in F_l[\text{Pr } S] \wedge d \in F_l[\text{Pr } T]}(R) \right| \right), \quad (22)$$

$$F_l \otimes (x \in X) = F_l' = \left( F_l[\text{Pr } S] \times X; F_l[\text{Pr } T]; \left| \sigma_{x \in X} \left( \sigma_{k \in F_l[\text{Pr } S]}(R) \right) \right|; \left| \sigma_{x \in X} \left( \sigma_{k \in F_l[\text{Pr } S] \wedge d \in F_l[\text{Pr } T]}(R) \right) \right| \right). \quad (23)$$

Тут  $\sigma_{\text{Pr}}(R)$  – операція селекції по предикату Pr над відношенням R.

Поки не вдалось розробити ефективний метод обчислення множин  $\sigma_{x \in X} \left( \sigma_{k \in F_l[\text{Pr } S]}(R) \right)$  та  $\sigma_{x \in X} \left( \sigma_{k \in F_l[\text{Pr } S] \wedge d \in F_l[\text{Pr } T]}(R) \right)$  по наявних  $\sigma_{k \in F_l[\text{Pr } S]}(R)$  та  $\sigma_{k \in F_l[\text{Pr } S] \wedge d \in F_l[\text{Pr } T]}(R)$  (їх можна зберігати для кожної аналізованої залежності в ході обчислень для спрощення побудови мультиатрибутних залежностей). Єдиним виходом залишається пряме виконання вказаних операцій селекції.

Отже, алгоритм пошуку ПЗ розширюється додатковим кроком – пошуком мультиатрибутних залежностей. Для цього залежності з кроку 1 алгоритму аналізуються на наявність кореляцій з іншими атрибутами відношення. Повний перебір можливих розширень ПЗ до мультиатрибутних буде мати складність  $O(z^2 \cdot m \cdot sz)$ , де  $z$  – кількість знайдених ПЗ,  $m$  – кількість атрибутів відношення,  $sz$  – середній розмір множини кортежів відношення R, на яких визначена ПЗ.

Для оптимізації цього процесу можна перебирати пари залежностей виду (6) –  $(F_1^A; F_2^A)$ , які мають значне перекриття множин значень результуючої частини предиката:

$$\frac{F_1^A[\text{Pr } T] \cap F_2^A[\text{Pr } T]}{F_1^A[\text{Pr } T] \cup F_2^A[\text{Pr } T]} \geq c. \quad (24)$$

Значення  $c$  визначається експертним шляхом, залежно від структури даних та типів залежностей, які слід шукати.

Пари залежностей, що задовольняють умову (24), комбінуються операціями об'єднання (якщо умовні частини предиката мають однакові атрибути) чи додавання атрибуту в умову предиката (якщо умовні частини предиката мають різні атрибути), а утворені залежності додаються у множину ПЗ для подальшого агрегування.

Використання умови (24) дозволяє значно скоротити перебір можливих мультиатрибутичних залежностей і забезпечити складність виконання другого кроку алгоритму  $O(z^2 \cdot \log(sz) + z^* \cdot sz)$ , де  $z^*$  – кількість пар ІПЗ, що задовольняють умову (24). Обчислення формули (24) може бути виконано з використанням структури даних binomial heap або Fibonacci heap зі складністю  $O(\log(sz))$ . Таким чином, якщо кількість початкових залежностей не надто велика (порядку декількох тисяч), досягається прийнятний час аналізу мультиатрибутичних залежностей на сучасних обчислювальних системах. Контролювати кількість одноатрибутичних залежностей можна вибором порогового значення імовірності  $p_0$ , а перебір мультиатрибутичних залежностей – зміною константи  $c$  з умови (24).

Таким чином, описаний у даній статті алгоритм дозволяє проводити ефективний аналіз однотипних даних на наявність мультиатрибутичних ІПЗ та ЧФ, сумарна складність якого  $O(m^2 \cdot n + z^2 \cdot \log(sz) + z^* \cdot sz)$ . При розумному виборі значень  $p_0$  та  $c$  на сучасних обчислювальних центрах він дає змогу аналізувати набори даних з тисячами атрибутів та мільйонами кортежів відношення.

Розглянемо продовження прикладу 1 для демонстрації пошуку мультиатрибутичних залежностей. Виберемо значення  $c = 0,7$ . Розглянемо одноатрибутичний ІЗ  $P(L=1 \rightarrow T \in \{2,3\}) = 1$  та  $P(i \in \{1,2,3\} \rightarrow T=3) = \frac{2}{3}$ ; за формулою (24) обчислюємо їх перекриття:

$$\frac{F_1^A[\text{Pr } T] \cap F_2^A[\text{Pr } T]}{F_1^A[\text{Pr } T] \cup F_2^A[\text{Pr } T]} = \frac{3}{4} \geq c.$$

Отже, обчислюємо мультиатрибутичну залежність, утворену додаванням атрибуту  $i$  на множині значень  $\{1,2,3\}$  в залежність  $P(L=1 \rightarrow T \in \{2,3\}) = 1$ :

$$P(L=1 \wedge i \in \{1,2,3\} \rightarrow T \in \{2,3\}) = \frac{|\sigma_{i \in \{1,2,3\}}(\sigma_{L=1}(R))|}{|\sigma_{i \in \{1,2,3\}}(\sigma_{L=1 \wedge T \in \{2,3\}}(R))|} = \frac{3}{3} = 1.$$

Таким чином, утворилась нова мультиатрибутична ІПЗ.

### 3. Висновки

Описаний у даній статті алгоритм дозволяє проводити аналіз великих обсягів даних, які можуть бути представлені у вигляді відношення (на сучасних обчислювальних центрах за кілька годин можна провести аналіз відношення з тисячами атрибутів та мільйонами записів). Алгоритм достатньо легко можна модифікувати для паралельного виконання на декількох обчислювальних машинах, що при потребі дозволить зменшити час аналізу до хвилин.

Система аналізу даних, яка використовуватиме вищеописаний алгоритм, дозволить значно розширити межі застосування методів аналізу даних у напрямку розміру аналізованих відношень, що, у свою чергу, дозволить застосовувати їх до предметних галузей, в яких накопичуються надзвичайно великі архівні бази даних, та виявити нові важливі залежності в цих галузях.

Основним недоліком розробленого алгоритму є те, що він не використовує вже наявних знань про специфіку предметної галузі і дозволяє виявити далеко не всі типи можливих залежностей (його завданням є пошук імовірнісних продукційних залежностей). Ще одним недоліком є недостатня ефективність пошуку мультиатрибутичних залежностей і для



оптимізації цього процесу доводиться вводити додатковий коефіцієнт, на якому базується евристичний алгоритм пошуку таких залежностей.

І все ж, не зважаючи на описані недоліки, алгоритм дозволяє виявляти дуже широкий клас залежностей у даних і може бути застосований у найрізноманітніших галузях науки і техніки.

## СПИСОК ЛІТЕРАТУРИ

1. Головний сайт департаменту патології, UT Southwestern Medical Center [Електронний ресурс]. – Режим доступу: <http://pathcuric1.swmed.edu/pathdb/classifi.html>.
2. Опис утиліти BiNGO, сайт університету Гент [Електронний ресурс]. – Режим доступу <http://www.psb.ugent.be/cbd/papers/BiNGO/Home.html>.
3. Офіційний сайт National Institute of Allergy and Infectious Diseases (NIAID), NIH [Електронний ресурс]. – Режим доступу <http://david.abcc.ncifcrf.gov/content.jsp?file=/ease/ease1.htm&type=1>.
4. Мейер Д. Теория реляционных баз данных / Мейер Д.; пер. с англ. – М.: Мир, 1987. – 610 с.

*Стаття надійшла до редакції 14.07.2011*