

Н. Шаховська

Національний університет “Львівська політехніка”,
кафедра інформаційних систем та мереж

МОДЕЛЬ КОНСОЛІДОВАНИХ ДАНИХ ТА ЇХ ОПРАЦЮВАННЯ ЗА УМОВ НЕВИЗНАЧЕНОСТІ

© Шаховська Н., 2013

Проаналізовано проблеми опрацювання розрізнених даних. Побудовано формальну модель сховища консолідованих даних та розроблено метод зменшення невизначеності даних.

Ключові слова: невизначеність даних, модель сховища, консолідовані дані.

Problems which arise up during work with separate sources with depositories information using and databases are analyzed. There is formalized model of consolidation datawarehouse and there is built the methods for uncertainty elimination .

Key words: uncertainty data, warehouse model, consolidated data.

Вступ

Інформаційні об'єкти, що описують певну предметну область, консолідовані дані та зв'язки між об'єктами, становитимуть простір даних. Однією зі задач, яка виникатиме у процесі консолідації, є невизначеність даних, що є результатом дублювання, неточності, відсутності, суперечливості даних. Також невизначеність виникатиме внаслідок встановлення неправильних зв'язків між об'єктами. Тому виникає задача зменшення невизначеності для підвищення якості даних.

Оскільки дані надходять з різних джерел, певна множина даних може бути відсутня у джерелах даних, а інша може перекриватися у різних інформаційних продуктах. Тому виникає проблема дублювання, відсутності, неповноти та нечіткості даних.

Невизначеність може виникати на рівні атрибута, кортежа та відношення (невизначеність у схемі опису). Поява невизначеності на рівні атрибута і кортежа у зв'язку з багатовимірністю відображення інформації призводить до поширення невизначеності на всі примірники опису певного концепту. Оскільки простір даних об'єднує мільйони даних про об'єкти проблемної області, то опрацювання невизначеності традиційними засобами (інтервальна математика, багатозначна логіка) стає неефективним через велику кількість операндів.

Отже, специфіка простору даних (наявність множини різнотипних джерел, дублювання даних, неоднозначність опису джерел даних) приводить до того, що невизначеність, яка у традиційних реляційних базах даних розглядалася у межах одного відношення і могла виникати на рівні атрибута, кортежа та на рівні відношення, в цьому випадку поширюється через сприйняття користувачем інформації на весь простір даних. Тому для опрацювання невизначеності у просторі даних необхідно використати якісно новий підхід, потреба застосування якого не виникала у реляційних базах даних та сховищах даних.

Аналіз літературних джерел і постановка задачі

Прокласифікуємо типи невизначеності за характером їх появи у просторі даних. Однією з перших робіт у цьому напрямі є робота Л. Заде [1]. Г. Цельмер [2] підкреслює, що невизначеність, будучи об'єктивною формою існування реального світу, обумовлена, з одного боку, об'єктивним існуванням випадковості як форми прояву необхідності, а з іншого – неповнотою кожного акту відображення реальних явищ у людській свідомості. Причому неповнота відображення принципово

непереборна через загальний зв'язок всіх об'єктів реального світу і нескінченність їх розвитку. Полягає невизначеність у різноманітті перетворення можливостей на дійсність, в існуванні множини (як правило, нескінченної кількості) станів, в яких об'єкт, що змінюється в динаміці, може перебувати в майбутній момент часу.

Ф. Найт під ситуацією невизначеності розуміє недостатню обізнаність та необхідність діяти, спираючись на думку, а не на знання [3]. В.В. Черкасов трактує невизначеність як постійну мінливість умов, швидко й гнучку переорієнтацію виробництва, дії конкурентів, зміну ринку і т.п. Невизначеність він називає найхарактернішою причиною ризику в управлінській діяльності. Він виділяє два класи невизначеності: “доброякісна” невизначеність, коли для невідомих чинників є статистичні або імовірнісні характеристики, і “погана” невизначеність, коли таких характеристик в принципі не можна отримати, причому існують методи означення обох видів невизначеності, що виникають в реальних завданнях (Вентцель, 1980).

У (Моїсеєв, 1975) наводиться така класифікація невизначеностей [4]:

- за ступенем невизначеності: імовірнісна, лінгвістична, інтервальна, повна невизначеність;
- за характером невизначеності: параметрична, структурна, ситуаційна;
- за використанням одержаної в ході керування інформації: переборна і невиважна.

У Дієва В.С. і Трухачева Р.І. [5] наведено детальнішу класифікацію невизначеностей у сучасних економічних системах.

У [6] визначено типи невизначеностей, природою яких є:

- 1) значення невідоме (відсутнє);
- 2) неповнота інформації;
- 3) нечіткість (стохастичність) – використання розподілу для встановлення істинності знань;
- 4) неточність (стосується числових даних);
- 5) недетермінованість процедур виведення рішень (випадковість);
- 6) ненадійність даних;
- 7) багатозначність інтерпретацій;
- 8) лінгвістична невизначеність: невизначеність значення слова, невизначеність змісту речення.

Для сховища даних вважалося, що невизначеність може виникати на рівні значення та на рівні кортежу [7]. Оскільки для простору даних притаманним є встановлення факту довіри до джерела, то рівні уведення невизначеності необхідно розширити.

Змінимо рівні уведення типів невизначеностей у просторі даних (рис. 1).

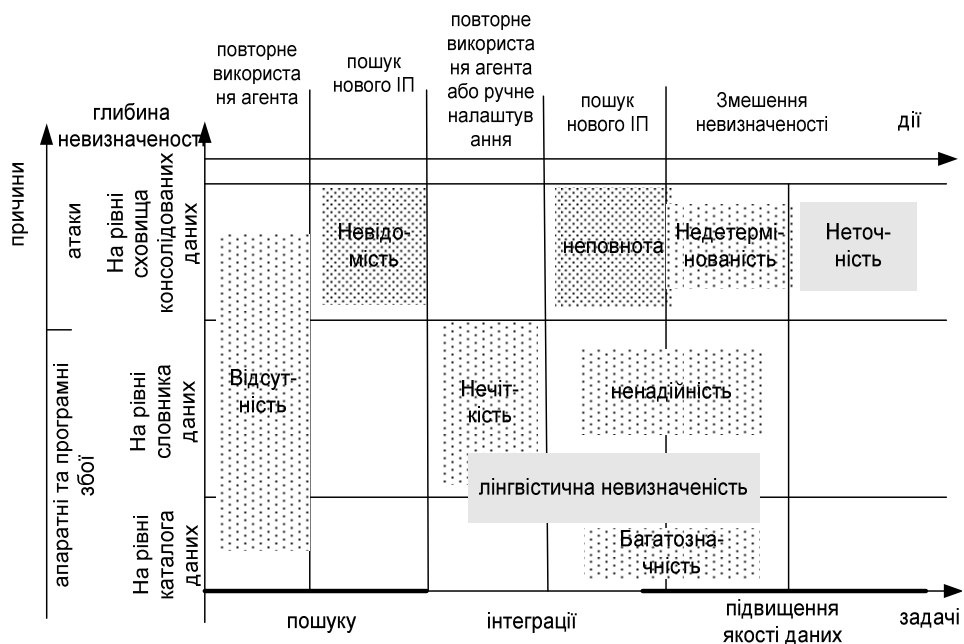


Рис. 1. Типи невизначеності у сховищі консолідованих даних простору даних та рівні їх введення

Невизначеності на рівні сховища даних виникають на основі атак – блокування даних у джерелі, приховуванням частини інформації тощо.

Невизначеності на рівні словника та каталогу даних виникають насамперед на основі програмних збоїв, а також через наявність атак на рівні джерел даних.

Розглянемо детальніше вказані типи невизначеностей та виявимо місця їх появи у сховищах та просторах даних [7]. Аналізуватимемо невизначеності, що виникають у результаті консолідації даних у *єдине джерело* (локальне чи віртуальне), а, отже, матимемо справу зі структурованими даними. Для подання *єдиного джерела* використовуватимемо *реляційну модель*.

Відсутність даних виникає внаслідок відсутності опису необхідної характеристики у каталозі даних та словнику. Відсутність може виникнути або через те, що необхідної характеристики не знайдено у інформаційних продуктах простору даних, або вона не включена до каталогу чи словника через недостатній рівень довіри. Для усунення цього виду невизначеності необхідно повторне використання агента, можливо зі зменшеним рівнем довіри до даних.

Невідомість даних зустрічається на рівні значення характеристики (атрибуту у реляційних базах даних) і означає, що значення притаманне об'єкту, але невідоме.

$$s = \{A, unk\},$$

де s – об'єкт, який описується кортежем характеристик консолідованих даних, unk – відсутнє значення, A – решту значень атрибутів характеристик кортежу консолідованих даних, $unk \cup A = s, unk \cap A = \emptyset$.

Подання цього типу невизначеності є ідентичним до подання у сховищах даних. Поява невідомості на рівні каталогу даних призводить до зашумлення всієї інформації, що отримується від джерела даних з невідомим атрибутом.

Неповнота є станом об'єкта, у якому є підмножина відсутніх значень характеристик. Якщо ця підмножина є порожня і ми говоримо про реляційне подання даних, то отримаємо традиційний кортеж. Відсутність інформації є також частковим випадком неповноти інформації, коли кількість невідомих значень атрибутів кортежу дорівнює 1. Неповнота може з'являтися як і у відношенні, у яке інтегруються дані, так і у словнику даних як результат збоїв роботи інтелектуального агента визначення структури джерела:

$$s = \{A, \{unk\}\}, |unk| < |A|.$$

Моделюється так само, як і в сховищах даних, але, на відміну від сховища даних, виникає і на рівні відношення (каталогу даних).

Невизначеності типів 3 – 8 класифікують як неоднозначність даних, що переважно виникають на рівні об'єкта або підмножини значень характеристик, із яких формується кортеж. Вони виникають як результат атак на рівні джерел даних (інформаційних продуктів).

Нечіткість виникає через неповне вивчення або неоднозначне відображення характеристик сутності. Моделюється за допомогою доповнення схеми відношення додатковим атрибутом (атрибутами), значення яких містять рівень впевненості в істинності підмножини значень неключових атрибутів.

$$s = \{A, unk_1, unk_2, \dots, unk_n\}, A \in K, A', 1 \leq n \leq |A|, unk_{attr} = P^{attr}(i, j), |A| \geq \{unk_1, unk_2, \dots, unk_n\}$$

де K – множина значень ключів, A' – підмножина значень неключових атрибутів. Рівень впевненості може позначатися за допомогою числової шкали, лінгвістичних оцінок, нечіткої величини тощо.

Неточність отримується внаслідок застосування математичних операцій над числовими даними (цього типу є також невизначеність, яка виникає внаслідок роботи з інтервальними величинами). Цей тип невизначеності моделюється за допомогою додаткового атрибута і може виникати через нечіткість в словнику даних.

На відміну від сховищ даних, у просторах даних виникає доволі часто у зв'язку з опрацюванням даних, що зберігались на різних платформах, використовувались для вирішення різного класу задач.

$$s = \{A, \{unk\}\}, \{unk\} \subset A, Design(A) \in \{unk\}.$$

Недетермінованість процедур виведення рішень (випадковість) виникає у випадку, коли необхідно зберігати проміжні або кінцеві результати процедур виведення або прийняття рішень, а також – у відношенні фактів на рівні значень агрегованих атрибутів. Моделюється за допомогою розширення схеми даних та виникає винятково у сховищі консолідованих даних:

$$s = s \cup \{unk\}, \{unk\} \notin \mathbf{A}, Design(s) \in \{unk\}.$$

Ненадійність є типом невизначеності, який вважається однією із характеристик об'єкта. Хоча сама природа цієї характеристики є невизначеною, у відношенні як її домен використовують традиційну числову шкалу та застосовують до її значень традиційні математичні операції. Виникає внаслідок визначення довіри звернення до джерела даних. Моделюється за допомогою доповнення схеми каталогу даних додатковим атрибутом. Значення цього атрибута змінюється у результаті роботи простору даних. Представляється як характеристика, обернена до значення довіри до джерела даних.

$$s = s \cup [unk_j], unk_j \notin \mathbf{A}, unk_j = \frac{1}{P(j)}.$$

Багатозначність інтерпретації є одним із джерел виникнення суперечностей. Такий тип невизначеності виникає найчастіше у каталозі даних через отримання інформації із різних джерел і неможливість визначення істинності даних. Для відображення цього типу невизначеності схему відношення доповнюють додатковим атрибутом, який містить ступінь впевненості у істинності даних кортежу. Від типу нечіткості відрізняється тим, що вводиться на рівні відношення.

Лінгвістична невизначеність пов'язана з використанням природної мови в інформаційних ресурсах (у текстових файлах та веб-ресурсах), які мають якісний характер, і може виникати внаслідок нерозуміння (незнання) значення слова або нерозуміння змісту речення. Такий тип невизначеності зустрічається у системах опрацювання текстової інформації (системи автоматизованого перекладу, системи для самонавчання тощо). У контексті просторів даних виникає внаслідок опрацювання напівструктурованої інформації (тексти, веб-сторінки тощо).

Розглянуті типи невизначеностей можуть накладатись або бути джерелом появи одна однієї.

Для задачі *зменшення невизначеностей* удосконалено метод, який використано для зменшення невизначеності у сховищах даних реляційного типу – усунення невизначеності на основі методу видобування знань. Невідоме значення атрибута розглядається як позначка класу, а сама задача усунення невизначеності трансформується у задачу віднесення до класу. Використання цього методу дає змогу усувати невизначеності типу “невідомий” та “неповний” на рівні значення атрибута та підмножини атрибутів. Проте, на відміну від сховищ даних, необхідно ще враховувати рівень довіри до джерела даних, тобто працювати з невизначеністю на рівні відношення.

Задача усунення (зменшення) невизначеностей – це побудова гомоморфного відображення множини даних, що зберігаються у сховищі консолідованих даних, у множину даних, що використовуються для підтримання прийняття рішень (рис. 2) з метою підвищення якості консолідованих даних та прийняття на їх основі ефективних керівних стратегічних рішень, враховуючи ймовірність появи атак. Атака – додавання до простору даних джерела, структура даних яких викликає багатозначність інтерпретацій у словнику синонімів *Dis*.

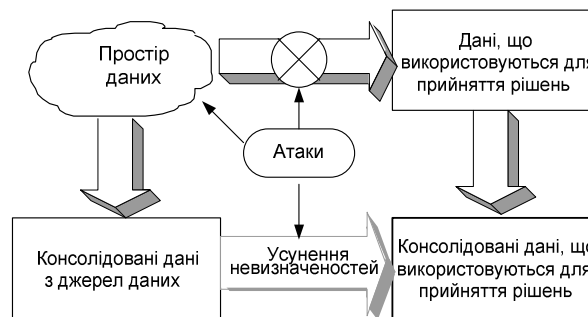


Рис. 2. Постановка задачі усунення невизначеності у сховищі консолідованих даних простору даних

Одним із методів моделювання неточних, нечітких та часткових даних є введення у каталог джерела додаткового атрибута, значення якого визначає ступінь довіри до невизначених даних.

Основний матеріал

1. Модель сховища консолідованих даних

Схема сховища консолідованих даних Cg' – скінченна множина імен атрибутів $\{A_1, A_2, \dots, A_n\}$, значення яких є чіткими; $\{A_unk_1, A_unk_2, A_unk_p\}$ з нечіткими або недетермінованими значеннями; множину імен атрибутів $\{Unk_1, Unk_2, \dots, Unk_m\}$, доменами яких є числові дані, що моделюють імовірнісні дані, значення функції приналежності нечітких множин, ступінь істинності багатозначної логіки, процентні відношення, коефіцієнти, різноманітні шкали або лінгвістичні оцінки; схему словника синонімів Dic та схему каталогу даних Cg :

$$Cg' = \langle \{C_1, C_2, \dots, C_n\}, \{C_unk_1, C_unk_2, C_unk_p\}, \{Unk_1, Unk_2, \dots, Unk_m\}, Dic, Cg \rangle, (3.1)$$

Невизначеними вважаються значення атрибутів множини C_unk , а рівень довіри до них зберігається у значеннях атрибутів множини Unk .

Кортеж консолідованих даних dc – інформаційний опис об'єкта t джерела даних S , поданий у вигляді множини (кортежу) значень характеристик (атрибутів), підмножина значень атрибутів якого містить дані про об'єкт, джерело даних та синонімічні назви об'єкта, причому ці дані можуть бути неповні, нечіткі чи недетерміновані дані. Тобто, об'єкт, який моделюється у джерелі даних цим кортежем, існує, але частина інформації про нього відсутня, нечітка, неповна, недетермінована тощо.

Значення атрибутів кортежу консолідованих даних поділимо на групи.

1. Чіткі (відомі) – значення первинного ключа, зовнішніх ключів (можуть бути відсутні). Позначимо їх через C .

2. Відсутні – фізично відсутня інформація. Позначимо їх через \perp .

3. Невизначені – для підмножин атрибутів введена множина атрибутів Unk , які вказують ступінь істинності значень цих атрибутів. За замовчуванням значенню атрибута Unk присвоюємо значення, яке означає найвищий ступінь істинності.

Зауважимо, що, у випадку стовідсоткової довіри до кожного значення кортежу, ми отримуємо традиційний реляційний кортеж та застосовуємо традиційні операції над ним.

Кортеж консолідованих даних dc – це множина значень характеристик об'єкта сутності:

$$dc = \langle C, C_unk, Unk, \{dic\}, \{cg\} \rangle,$$

де C – підмножина значень атрибутів із чіткими значеннями, $C = t_{rel} \cup t_{dw} \cup t_{text}$, C_unk – підмножина значень атрибутів із нечіткими та недетермінованими значеннями, Unk – підмножина значень атрибутів із ступенями істинності значень атрибутів C_unk , $\{dic\}$ – множина значень словника даних, $\{cg\}$ – множина значень каталогу даних.

Сховище консолідованих даних cg' – множина відношень зі схемою Cg' та множиною кортежів консолідованих даних dc

Модель сховища консолідованих даних містить дані з усіх типів джерел простору даних.

2. Розроблення операцій над консолідованими даними моделі сховища

Оскільки сховище консолідованих даних є розширенням сховища даних, побудованого на основі реляційної моделі, то далі удосконалимо операції для роботи з ним.

Для опрацювання та аналізу невизначеностей за допомогою запиту в реляційних операторах слід здійснювати селекцію кортежів за значеннями множини атрибутів Unk . У сховищі даних аналогічною до неї є операція зрізу. Нехай r та s – відношення зі схемою R , r' та s' – відношення зі схемою $R \cup Unk \cup Dic \cup Cg$. Тоді $r \cap s$, $r \cup s$ і $r - s$ є відношеннями зі схемою R , а $r' \cap s'$, $r' \cup s'$ і $r' - s'$ – відношеннями зі схемою $R \cup Unk \cup Dic \cup Cg$.

Враховуючи ймовірність атак (невизначеність типу “багатозначність”), вибираємо ті джерела даних, рівень довіри до яких вищий за аналогічні:

$$r' = r \cup \sigma_{\max(P(\pi(Cg)))}(Dic) \cup Cg.$$

Доповнення до відношення r' працюватиме коректно у разі присвоєння всім значенням атрибута Unk найнижчого ступеня довіри (апріорі вважається, яка ця інформація, що заноситься у відношення є правдивою та повною, а про решту інформації нам нічого не відомо). Такий метод подання ступеня істинності за замовчуванням обираємо за принципом замкненості.

Оператор зрізу передбачає аналіз нечіткого значення за множиною значень атрибутів Unk .

$$slice : \sigma^{cons}_{\left(\frac{Unk \Theta unk}{\cup \sigma_C(Dic)} \cup \frac{C_unk \Theta c_unk}{\cup \sigma_C(Cg)}\right)}(cg') = \left\{ t \in dc \mid t(Unk) \Theta unk, t(C_unk) \Theta c_unk, meta_{Unk, C_unk} = 1, \right. \\ \left. \sigma_C(Dic) \text{ Is Not NULL}, \sigma_C(Cg) \text{ Is Not NULL}, unk = P(cg') \right\},$$

де Θ – множина символів (знаків) бінарних відношень над парами значень доменів. Вважається, що до кожного атрибута C_unk застосовуються операції порівняння. Як правило, будуть вживатися лише такі знаки порівняння над одним доменом: =, ≠, <, ≤, ≥, >.

Твердження: Удосконалений оператор зрізу, як і оператор вибірки, зберігає властивості комутативності та дистрибутивності відносно булевих операцій.

Доведення

Нехай $r'(R')$ – відношення, $R' \leftarrow R \cup Unk \cup Dic \cup Cg$, A і B – атрибути в R' , і нехай $a \in dom(A)$, $b \in dom(B)$. Тоді має місце рівність: $\sigma^{cons}_{A=a}(\sigma_{B=b}(r')) = \sigma_{cons_{B=b}}(\sigma_{A=a}(r'))$.

Удосконалений оператор зрізу дистрибутивний відносно бінарних булевих операцій:

$$\sigma^{cons}_{A=a}(r' \gamma s') = \sigma^{cons}_{A=a}(r') \gamma \sigma_{A=a}(s'),$$

де $\gamma = \cap, \cup$ або $-$, а r' і s' – відношення над однією і тією ж схемою.

Аналогом операції згортання у сховищі даних, побудованому на основі реляційної моделі, є *операція проєкції*. Здійснюючи проєкцію відношення з кортежами консолідованих даних, слід відслідковувати зв'язок підмножини атрибутів Unk із підмножиною атрибутів C_unk , а також перевіряти, чи для назви атрибута C_unk є синонім у словнику синонімів Dic . Тому удосконалений оператор згортання подано так:

$$drill\text{-}down : \pi_X^{cons}(cg') = \Pi F \left(\begin{array}{l} \neg ISNULL(\sigma_{Cg=R \cup C_unk=X}(c_unk)); \pi_{X \cup \pi_{Unk}(\sigma_{Cg=meta(C_unk, Unk)=1}(c_unk))}(dc); \\ \Pi F(\sigma_{C \cup C_Unk=X}(Dic); \pi_{\sigma_{C \cup C_Unk=X}(Dic)}(r); \pi_X(dc) \end{array} \right),$$

де $\Pi F(\text{умова}; \text{дія } 1; \text{дія } 2)$ – операція, введена у стандарті SQL 92. У разі виконання умови виконується дія 1, інакше дія 2; $ISNULL(r)$ – логічний оператор, результатом якого є істина, якщо відношення-операнд r не містить кортежів, та хибя – у іншому випадку. Також здійснюється пошук синоніма атрибута у словнику синонімів Dic ($\sigma_{C \cup C_Unk=X}(Dic)$) та заміна за потреби ($\pi_{\sigma_{C \cup C_Unk=X}(Dic)}(r)$).

Твердження: Удосконалений оператор згортання зберігає властивості традиційного оператора проєкції.

Доведення: Якщо $X_1 \subseteq X_2 \subseteq \dots \subseteq X_m \subseteq R'$, то $\pi_{X_1}^{cons}(\pi_{X_2}^{cons}(\dots(\pi_{X_m}^{cons}(cg'))\dots)) = \pi_{X_1}^{cons}(cg')$.

Оператор з'єднання використовується для зв'язування відношення фактів та відношень вимірів у сховищі консолідованих даних, оскільки воно будується на основі реляційної моделі.

Традиційний оператор з'єднання не може використовуватися для сховищ та просторів даних з консолідованими даними, оскільки для статистичного аналізу необхідне з'єднання відношення фактів з відношеннями вимірів, а за наявності непорожньої підмножини атрибутів Unk у відношеннях фактів та вимірів таке з'єднання буде некоректним. Також на оператор з'єднання впливає той факт, що виникає необхідність з'єднання не лише за тими атрибутами, що вказані як вхідні параметри, але й перевіряти наявність синонімів у словнику синонімів Dic .

Для удосконалення оператора з'єднання слід розглянути випадки, коли відношення є повністю з'єднувальними або не повністю з'єднувальними. Для повністю з'єднувальних відношень введення множини атрибутів Unk не впливає на операцію з'єднання. Якщо значення множини атрибутів Unk містять міру невизначеності зовнішнього ключа відношення, з яким відбувається з'єднання, то ця міра невизначеності переноситься на всі решту значень атрибутів цього відно-

шення. У випадку неповної з'єднувальності значення атрибута Unk для кортежів підлеглої таблиці, які не потрапляють у відношення, вважатимуться такими, що мають найвищий ступінь довіри.

$$across : r \triangleright \triangleleft cg' = \text{PF}(\sigma_{\sigma_{C \cup C_{Unk=X}}(Dic)}; \pi_{\sigma_{C \cup C_{Unk=X}}(Dic)}(r \triangleright \triangleleft cg'); \pi_{(R,B,NVL(Unk, \min))}(r \triangleright \triangleleft cg')),$$

де r – традиційне відношення, cg' – відношення з консолідованими даними, R – множина атрибутів відношення r , S – множина атрибутів відношення cg' , не включаючи підмножини атрибутів Unk ($Cg' = Cg \cup Unk$), B – множина тих атрибутів з S , яких нема у відношенні r ($B \subset Cg$, $B \not\subset Cg \cap R$), \min – значення, яке означає найнижчий ступінь довіри, $NVL(Unk, \min)$ – операція, яка присвоює \min усім значенням Unk для нез'єднувальних кортежів відношення cg' , $\triangleright \triangleleft$ – ліве з'єднання.

Спочатку перевіряється, чи необхідно здійснювати з'єднання не за заданими атрибутами, а за синонімами ($\sigma_{\sigma_{C \cup C_{Unk=X}}(Dic)}$). Якщо ні, то виконується операція лівого з'єднання для відношень з схемами S' і R та проєкція за атрибутами-синонімами. В іншому випадку виконується операція лівого з'єднання за спільними атрибутами, потім над отриманим з попередньої операції відношенням здійснюється операція проєкції, за якою утвореним у результаті з'єднання порожнім значенням підмножини атрибутів Unk присвоюється значення \min .

Слід зазначити, що коли словник синонімів порожній ($Dic = \emptyset$) і ймовірність звернення до джерел даних загалом та до їх характеристик дорівнює одиниці ($Unk = 1$), то отримаємо традиційне реляційне з'єднання.

Твердження: Удосконалений оператор з'єднання комутативний та асоціативний.

Доведення

Для даних відношень q' , r' і s'

$$(q' \triangleright \triangleleft r') \triangleright \triangleleft s' = q' \triangleright \triangleleft (r' \triangleright \triangleleft s').$$

Введемо позначення для деяких багаторазових з'єднань. Нехай $s'_1(S'_1), s'_2(S'_2), \dots, s'_m(S'_m)$ – відношення, $R' = S'_1 \cup S'_2 \cup \dots \cup S'_m$ і S' – послідовність s'_1, s'_2, \dots, s'_m . Далі, нехай t_1, t_2, \dots, t_m – послідовність кортежів, в якій $t_i \in s'_i, 1 \leq i \leq m$. Кортежі з'єднувальні на S' , якщо існує кортеж $t \in R'$, такий, що $t_i = t(S'_i), 1 \leq i \leq m$. Кортеж t є результатом з'єднання кортежів $t_1, t_2, \dots, t_m \in S'$.

3. Зменшення невизначеності консолідованих даних

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, визначена 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 мають не чітко детермінований характер; називатимемо їх імовірнісними продукційними залежностями (ІПЗ).

Імовірнісна продукційна залежність (ІПЗ) – 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.

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

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 ПЗ до кількості об'єктів в селекції:

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

Класифікаційним правилом назвемо імовірнісну продукційну залежність між підмножинами атрибутів X та Y у сховищі консолідованих даних cg' , яка зустрічається у тестовому наборі cg' зі ступенем відповідності (довіри) s , де $(X = x) \rightarrow (Y = y)$.

Будується класифікаційне правило на основі навчального набору даних у cg' , де значення позначок класу (значення підмножини атрибутів Y) відомі. Класифікаційне правило будується для схеми Cg' і тому не залежатиме від нових кортежів, що надходять у відношення сховища консолідованих даних (незалежність від тестового набору).

Позначкою класу назвемо лінгвістичну змінну або звичну характеристику об'єктів, яка є значеннями підмножини атрибутів Y і позначає об'єкти зі спільними (подібними зі ступенем s) значеннями підмножини атрибутів X . Домени атрибутів, що належать до підмножини Y , $y \in \text{dom}(Y) = \pi_Y(Cg')$, обов'язково повинні містити скінченну та наперед відому множину значень.

Позначки класу обираються з наперед відомої множини значень (у межах досліджуваної області є фіксованими), а віднесення до класу об'єктів, інформація про які щойно надійшла у сховище консолідованих даних, здійснюється на основі класифікаційних правил.

Додавання позначок здійснюватиметься автоматично, оскільки надходження нових джерел даних у простір даних – також динамічне.

Обчислення достовірності виконання такої залежності ґрунтується на можливості розкладу такої залежності на складові ПЗ:

$$P(s \in S \rightarrow t \in T) = \sum_{t_i \in T} P(s \in S \rightarrow t = t_i) = \sum_{t_i \in T} \frac{\sum_j |s = s_j \wedge t = t_i|}{\sum_j |s = s_j|} \quad (4)$$

Як і у випадку з F-залежностями (функціональними залежностями), множину класифікаційних правил, що існують у заданому відношенні, можна представити деякою їх підмножиною, з якої за допомогою правил виведення можна отримати усі класифікаційні правила даного відношення. Оскільки класифікаційні правила є розширенням F-залежностей, то варто розглянути трансформації аксіом для функціональних залежностей для класифікаційних правил.

Рефлексивність. $P(s \in S \rightarrow s \in S) = 1$ для будь-якого відношення $r(R)$.

$$\text{Доведення: } P(s \in S \rightarrow s \in S) = \frac{|\sigma_{s \in S \wedge s \in S}|}{|\sigma_{s \in S}|} = \frac{|\sigma_{s \in S}|}{|\sigma_{s \in S}|} = 1$$

Поповнення. Якщо $P(s \in S \rightarrow t \in T) = p$, то $P(s \in S \wedge w \in D(W) \rightarrow t \in T) = p$.

$$\text{Доведення: } P(s \in S \wedge w \in D(W) \rightarrow t \in T) = \frac{|\sigma_{s \in S \wedge w \in D(W) \wedge t \in T}(R)|}{|\sigma_{s \in S \wedge w \in D(W)}(R)|} = |\forall x \in r: q = \pi_{W=W}(x) \in D(W) \Rightarrow w \in D(W)| =$$

$$= \frac{|\sigma_{s \in S \wedge t \in T}(R)|}{|\sigma_{s \in S}(R)|} = P(s \in S \rightarrow t \in T) = p$$

Адитивність. Якщо $P(s \in S \rightarrow t \in T) = p$ і $P(s \in S \rightarrow w \in W) = 1$, то $P(s \in S \rightarrow t \in T \wedge w \in W) = p$.

$$\text{Доведення: } P(s \in S \rightarrow t \in T \wedge w \in W) = \frac{|\sigma_{s \in S \wedge t \in T \wedge w \in W}|}{|\sigma_{s \in S}|} = |s \in S \rightarrow w \in W| = \frac{|\sigma_{s \in S \wedge t \in T}|}{|\sigma_{s \in S}|} = P(s \in S \rightarrow t \in T) = p$$

Усунення невизначеностей, які зустрічаються серед значень атрибута Y відношення r , є класифікуванням із використанням модифікованого алгоритму *chase*.

Суть методу:

- 1) пошук кортежів, у яких є однакові значення по множині атрибутів X ;
- 2) пошук кортежів, у яких є однакові значення по множині синонімів атрибутів X ;
- 3) розрахунок рівня довіри до джерела кортежів, отриманих на кроках (1) та (2);
- 4) розрахунок рівня довіри до атрибутів джерел кортежів, отриманих на кроках (1) та (2).
- 5) визначення кортежів з найбільшими рівнями довіри.

Щоб мати можливість класифікувати об'єкти, необхідно побудувати функції класифікації. Взагалі, у просторі даних може зберігатися інформація про декілька типів класів, і для кожного типу класу є своя підмножина функцій. Одна й та ж функція може застосовуватись для визначення кількох типів класів.

Функціями класифікації називатимемо модифіковані функціональні залежності, які виконуються для певної підмножини кортежів відношення сховища консолідованих даних.

Алгоритм віднесення до класу:

1. **Якщо** $\sigma(cg') = \{dc_1(X_1) \downarrow, \dots, dc_1(X_n) \downarrow\} \mathbf{i} \{dc_2(X_1) \downarrow, \dots, dc_2(X_n) \downarrow\}$
 $\mathbf{i} \{dc_1(X_1) \downarrow, \dots, dc_1(X_n) \downarrow = dc_2(X_1) \downarrow, \dots, dc_2(X_n) \downarrow\}$
 $\mathbf{i} \{dc_1(Y) \downarrow\} \mathbf{i} \{dc_2(Y) = \perp\} \mathbf{i}$ **Якщо** $\sigma_{X_1}(Dic) = \emptyset$
то замінюємо \perp на $dc_1(Y) \mathbf{i} dc_1(P) = dc_1(P) / (\sum_i m_{1i} / n)$.
2. **Якщо** $\{dc_1(X_1) \downarrow, \dots, dc_1(X_n) \downarrow\}$
 $\mathbf{i} \{в\ dc_2\ m\ з\ n\ значень\ атрибутів - \downarrow, n - m\ значень\ атрибутів - \perp, m \leq n\}$
 $\mathbf{i} \{P \geq 1 - m/n\} \mathbf{i} \{по\ визначених\ значеннях\ dc_1(X^m) \downarrow = dc_2(X^m) \downarrow\}$
 $\mathbf{i} \{dc_1(Y) \downarrow\} \mathbf{i} \{dc_2(Y) = \perp\},$
то замінюємо \perp у r на $dc_1(Y) \mathbf{i} dc_2(P) = dc_2(P) / (\sum_i m_{2i} / n)$.
3. **Якщо** $\{в\ dc_i\ m_i\ з\ n\ значень\ атрибутів - \downarrow, m_i \leq n\}$
 $\mathbf{i} \{в\ dc_j\ m_j\ з\ n\ значень\ атрибутів - \downarrow, m_j \leq n\}$
 $\mathbf{i} \{по\ визначених\ значеннях\ dc_i(X^m) \downarrow = dc_2(X^m) \downarrow\}$
 $\mathbf{i} \{по\ визначених\ значеннях\ dc_j(X^m) \downarrow = dc_2(X^m) \downarrow\}$
 $\mathbf{i} \{m_i/n \leq m_j/n\} \mathbf{i} \{P \geq 1 - m_i/n\}$
 $\mathbf{i} \{dc_i(Y) \downarrow\} \mathbf{i} \{dc_j(Y) \downarrow\} \mathbf{i} \{dc_2(Y) = \perp\},$
то замінюємо \perp на $dc_j(Y) \mathbf{i} dc_2(P) = dc_2(P) / (\sum_i m_{2i} / n)$.

Приклад

Нехай маємо базу даних науково-дослідної частини та наукові звіти двох кафедр. Зрозуміло, що між вказаними джерелами існує залежність: потрапляння даних у базу даних науково-дослідної частини через наукові звіти кафедр. Нехай у сховище консолідованих даних у результаті інтеграції потрапило 2 кортежі з таким вмістом:

ID	Author	Title	Publisher	Source	Trust
1	Автор1, Автор2	Назва1	Журнал1	DB1	0,7
2	Автор1, Автор2	Назва1	Журнал2	Text32	0,4

Отримали невизначеність виду “багатозначність інтерпретацій” (відмінності значення атрибута **Publisher**).

У каталозі даних вказано, що *Text32* є джерелом даних для *DB1*. Тоді, не дивлячись, що рівень довіри до першого кортежу є вищим, ніж до другого, у результуючу вибірку потрапить другий кортеж.

Висновки

1. Побудовано модель сховища консолідованих даних, яка є розширенням моделі відношення з невизначеністю.
2. Удосконалено операції над відношенням з невизначеністю з метою їх застосування до сховища консолідованих даних, що дозволило реалізувати унарні операції алгебричної системи “сигнатури простір даних”.
3. Розроблено метод зменшення невизначеності даних, що розміщені у сховищі консолідованих даних як основу для подальшого оцінювання якості консолідованих даних.
4. Розглянуті методи доцільно застосовувати також і для прийняття рішень, так як це забезпечує пошук прихованих залежностей між характеристиками сховища консолідованих даних. Такі залежності доцільно враховувати під час прийняття рішень на основі консолідованих даних.

1. Заде Л. Понятие лингвистической переменной и его применение к принятию приближенных решений / Заде Л. – М.: Мир, 1976. – 166 с. 2. Цельмер Г. Учет риска при принятии управленческих решений / Г. Цельмер // Проблемы МСНТИ. – № 3. – С. 94–105, 1980. 3. Найт Ф.Х. Риск, неопределенность и прибыль / Ф.Х. Найт. – М.: Дело, 2003. – 358 с. 4. Моисеев Н.Н. Элементы теории оптимальных систем / Н.Н. Моисеев. – М.: Наука, 1975. – 528 с. 5. Трухачев Р.И. Модели принятия решений в условиях неопределенности / Р.И. Трухачев. – М.: Наука, 1981. – 151 с. 6. Згуровський М.З. Основи системного аналізу / Згуровський М.З., Панкратова Н.Д. – К.: Видавнича група ВНУ, 2007. – 544 с. 7. Шаховська Н.Б. Моделювання невизначеностей у сховищах даних реляційного типу. – Львів, 2007. – автореф. дис. ... канд. техн. наук.