

ФОРМАЛЬНА СЕМАНТИКА АГРЕГАТНИХ ОПЕРАЦІЙ МУЛЬТИМНОЖИННОЇ ТАБЛИЧНОЇ АЛГЕБРИ

І.М. Глушко

Ніжинський державний університет імені Миколи Гоголя,
16600, Ніжин, вул. Кропив'янського 2,
glushkoim@gmail.com

Розглянуто мультимножинну табличну алгебру. Сигнатура мультимножинної табличної алгебри поповнена агрегатними операціями. Задано формальну математичну семантику цих операцій та наведено приклади їх застосування.

Multiset table algebra is considered. The signature of multiset table algebra is filled up with aggregate operations. A formal mathematical semantics of these operations is defined.

Вступ

Реляційна модель даних на сьогодні широко використовується як у наукових дослідженнях в базах даних, так і на практиці. У формальному визначенні, запропонованому Е. Коддом [1], дана модель базується на множинах кортежів, тобто не дозволяє дублювати кортежів у відношенні. Багато мов, орієнтованих на роботу з базами даних, вимагають реляційну модель даних з мультимножинною семантикою (multi-set semantics) тому, що, по-перше, відношення, які дозволяють дублювати, корисні в багатьох прикладних областях, де об'єкти-дублікати можуть існувати; по-друге, в реляційній моделі даних видалення дублікатів після виконання операцій проєкції та об'єднання передбачає злиття однакових елементів або здійснення інших трудомістких дій.

Питанню використання мультимножин в базах даних приділяли увагу Paul W.P.J. Grefen та Rolf A. de By [2], G. Lamperti, M. Melchiori, M. Zanella [3], Г. Гарсія-Моліна, Дж. Ульман, Дж. Уйдом [4], A. Silbeschatz, H. Korth, S. Sudarshan [5], Д.Б. Буй, С.А. Поляков, Ю.Й. Брона, В.Н. Редько [6]. Разом з тим, в жодній із зазначених робіт не приділяється увага агрегатним операціям над таблицями мультимножинної табличної алгебри.

Мультимножини: основні поняття

Наведемо основні поняття мультимножин в термінах роботи [6]. Зафіксуємо деяку множину U . Під мультимножиною α з основою U будемо розуміти відображення вигляду $\alpha : U \rightarrow N$, де $N = \{1, 2, \dots\}$ – множина натуральних чисел.

Нехай D – універсум елементів основ мультимножин, тоді булеан $P(D)$ – універсум основ мультимножин. Під характеристичною функцією мультимножини α розуміємо функцію вигляду $\chi_\alpha : D \rightarrow Z_+$, значення якої задається наступною кусковою схемою:

$$\chi_\alpha(d) = \begin{cases} \alpha(d), & \text{якщо } d \in \text{dom } \alpha, \\ 0, & \text{інакше;} \end{cases}$$

для всіх $d \in D$, де Z_+ – множина цілих невід'ємних чисел.

Мультимножина називається порожньою і позначається як \emptyset_m , якщо її основа – порожня множина.

Мультимножини, областю значень яких є порожня множина або одноелементна множина вигляду $\{1\}$, називаються 1-мультимножинами. Дані мультимножини є аналогами звичайних множин.

Домовимося мультимножину α з основою $\{d_1, \dots, d_k\}$ записувати як $\{d_1^{n_1}, \dots, d_k^{n_k}\}$, де n_i – кількість дублікатів (екземплярів) елемента d_i у мультимножині α , $i = 1, \dots, k$.

Під рангом скінченної мультимножини α розуміємо суму дублікатів елементів її основи $\|\alpha\| = \sum_{d \in U} \chi_\alpha(d)$, де U – основа мультимножини α . Зрозуміло, що $\|\emptyset_m\| = 0$.

Мультимножина β включається у мультимножину α ($\beta \preceq \alpha$), якщо:

$$\beta \preceq \alpha \Leftrightarrow U_\beta \subseteq U_\alpha \ \& \ \forall d (d \in U_\beta \Rightarrow \beta(d) \leq \alpha(d)).$$

Якщо $\beta \preceq \alpha$, то мультимножина β називається підмультимножиною мультимножини α , а мультимножина α – надмультимножиною мультимножини β . Очевидно, що порожня мультимножина є підмультимножиною будь-якої мультимножини ($\emptyset_m \preceq \alpha$, $\forall \alpha$), а будь-яка мультимножина є своєю ж підмультимножиною ($\alpha \preceq \alpha$, $\forall \alpha$).

Мультимножинна таблична алгебра

Розглянемо дві множини: A – множину атрибутів і D – універсальний домен. Довільну скінченну множину атрибутів $R \subseteq A$ назовемо схемою.

Рядком схеми R називається іменна множина на парі R, D , проекція якої за першою компонентою рівна R (тобто по суті розглядається функція вигляду $s : R \rightarrow D$). Множину всіх рядків схеми R позначимо $S(R)$, а множину всіх рядків – S . Для позначення відсутніх значень у таблиці використовується особливий елемент універсального домену $NULL$. Позначимо s_R^{NULL} – константний рядок схеми R , тобто $s_R^{NULL} : R \rightarrow \{NULL\}$.

Задамо поняття таблиці як пари $\langle \psi, R \rangle$, де перша компонента ψ – це довільна мультимножина, зокрема, нескінченна, а друга компонента R – схема таблиці.

Таким чином, кожній таблиці приписується певна схема. Множину всіх таблиць схеми R позначимо $\Psi(R)$, а множину всіх таблиць $\Psi = \bigcup_{R \subseteq A} \Psi(R)$.

Позначимо $Occ(s, \psi)$ – кількість дублікатів (екземплярів) рядка s у мультимножині ψ . Домовимося мультимножину ψ записувати як $\{s_1^{n_1}, \dots, s_k^{n_k}, \dots\}$, де $n_i = Occ(s_i, \psi)$, $i = 1, 2, \dots$, а $\Theta(\psi) = \{s_1, \dots, s_k, \dots\}$ – основа мультимножин ψ .

Під мультимножинною табличною алгеброю розуміємо алгебру $\langle \Psi, \Omega_{P, \Xi} \rangle$, де Ψ – множина всіх таблиць, $\Omega_{P, \Xi} = \left\{ \bigcup_{All}^{\Psi, R}, \bigcap_{All}^{\Psi, R}, \setminus_{All}^{\Psi, R}, \sigma_{p, R}, \pi_{X, R}, \otimes_{R_1, R_2}, Rt_{\xi, R}, \sim_R \right\}_{\substack{p \in P, \xi \in \Xi \\ X, R, R_1, R_2 \subseteq A}}$ – сигнатура, P, Ξ – множини параметрів.

Задамо операції. Під об'єднанням \bigcup_{All}^R (перетином \bigcap_{All}^R , різницею \setminus_{All}^R) таблиць схеми R розуміється бінарна (параметрична) операція, отримана обмеженням однойменних операцій $\bigcup_{All}, \bigcap_{All}, \setminus_{All}$ над мультимножинами на множину всіх таблиць схеми R .

Розглянемо кожну операцію окремо. Позначимо через $\Theta(\psi_1)$ і $\Theta(\psi_2)$ – основи мультимножин ψ_1 та ψ_2 відповідно. Таким чином,

$$\bigcup_{All}^R : \Psi(R) \times \Psi(R) \rightarrow \Psi(R), \quad \langle \psi_1, R \rangle \bigcup_{All}^R \langle \psi_2, R \rangle = \langle \psi_1 \bigcup_{All} \psi_2, R \rangle.$$

Основа мультимножини $\psi_1 \bigcup_{All} \psi_2$ дорівнює об'єднанню основ мультимножин таблиць-аргументів:

$$\Theta(\psi_1 \bigcup_{All} \psi_2) = \Theta(\psi_1) \bigcup \Theta(\psi_2).$$

Дублікати рядків, які з'явилися після виконання операції, не вилучаються. Кількість дублікатів кожного рядка визначається за формулою:

$$Occ(s, \psi_1 \bigcup_{All} \psi_2) = \begin{cases} Occ(s, \psi_1), & \text{якщо } s \in \Theta(\psi_1) \setminus \Theta(\psi_2), \\ Occ(s, \psi_2), & \text{якщо } s \in \Theta(\psi_2) \setminus \Theta(\psi_1), \\ Occ(s, \psi_1) + Occ(s, \psi_2), & \text{якщо } s \in \Theta(\psi_1) \cap \Theta(\psi_2); \end{cases}$$

де $s \in \Theta(\psi_1) \bigcup \Theta(\psi_2)$.

$$\bigcap_{All}^R : \Psi(R) \times \Psi(R) \rightarrow \Psi(R), \quad \text{причому } \langle \psi_1, R \rangle \bigcap_{All}^R \langle \psi_2, R \rangle = \langle \psi_1 \bigcap_{All} \psi_2, R \rangle.$$

Основа мультимножини $\psi_1 \bigcap_{All} \psi_2$ дорівнює перетину основ мультимножин таблиць-аргументів:

$$\Theta(\psi_1 \bigcap_{All} \psi_2) = \Theta(\psi_1) \cap \Theta(\psi_2),$$

а кількість дублікатів рядка визначається як

$$Occ(s, \psi_1 \bigcap_{All} \psi_2) = \min(Occ(s, \psi_1), Occ(s, \psi_2)), \quad \text{де } s \in \Theta(\psi_1) \cap \Theta(\psi_2).$$

$$\setminus_{All}^R : \Psi(R) \times \Psi(R) \rightarrow \Psi(R), \quad \text{причому } \langle \psi_1, R \rangle \setminus_{All}^R \langle \psi_2, R \rangle = \langle \psi_1 \setminus_{All} \psi_2, R \rangle, \quad \text{де } \langle \psi_1, R \rangle, \langle \psi_2, R \rangle \in \Psi(R).$$

Основа мультимножини $\psi_1 \setminus_{All} \psi_2$ визначається як $\Theta(\psi_1 \setminus_{All} \psi_2) = \Theta(\psi_1) \setminus \Theta(\psi_2) \bigcup C^>(\psi_1, \psi_2)$, де $C^>(\psi_1, \psi_2) = \{s \mid s \in \Theta(\psi_1) \cap \Theta(\psi_2) \wedge Occ(s, \psi_1) > Occ(s, \psi_2)\}$.

Кількість дублікатів знаходиться так:

$$Occ(s, \psi_1 \setminus_{All} \psi_2) = \begin{cases} Occ(s, \psi_1), & \text{якщо } s \in \Theta(\psi_1) \setminus \Theta(\psi_2), \\ Occ(s, \psi_1) - Occ(s, \psi_2), & \text{якщо } s \in C^>(\psi_1, \psi_2), \end{cases}$$

де $s \in (\Theta(\psi_1) \setminus \Theta(\psi_2)) \cup C^>(\psi_1, \psi_2)$.

Нехай $p : S \xrightarrow{\sim} \{true, false\}$ – частковий предикат на множині рядків. Під селекцією за предикатом p таблиць схеми R розуміється унарна часткова параметрична операція $\sigma_{p,R}$, яка таблиці зіставляє її підтаблицю, що містить рядки, на яких предикат p істинний (крім того, предикат p має бути визначеним на всіх рядках таблиці-аргументу).

Отже,

$$\sigma_{p,R} : \Psi(R) \xrightarrow{\sim} \Psi(R), \text{ dom } \sigma_{p,R} = \{\langle \psi, R \rangle \mid \Theta(\psi) \subseteq \text{dom } p\}, \sigma_{p,R}(\langle \psi, R \rangle) = \langle \psi', R \rangle,$$

де $\langle \psi, R \rangle \in \Psi(R)$.

Областю означеності операції селекції є таблиці $\langle \psi, R \rangle \in \Psi(R)$, які містять рядки $s \in S(R)$, на яких предикат-параметр селекції p визначений. Надалі розглядаємо тільки предикати-параметри вигляду $p(s) = true \Leftrightarrow p(s(A_1), \dots, s(A_m)) = true$, де p – предикат на універсальному домені, а атрибути A_1, \dots, A_m належать схемі рядка s .

Основа мультимножини результуючої таблиці визначається як: $\Theta(\psi') = \{s \mid s \in \Theta(\psi) \wedge p(s) \simeq true\}$, де \simeq – узагальнена рівність (тобто обидві частини або одночасно не визначені або одночасно визначені і рівні).

В залежності від значення предиката на рядку s усі дублікати цього рядка або входять до мультимножини отриманої таблиці, або ні: $Occ(s, \psi') = Occ(s, \psi)$, де $s \in \Theta(\psi')$. Зазначимо, що операція селекції не породжує нові дублікати рядка.

Нехай $X \subseteq A$ – скінченна множина атрибутів. Під проекцією за множиною атрибутів X таблиць схеми R розуміється унарна параметрична операція $\pi_{X,R}$, значеннями якої є таблиці схеми $R \cap X$, що складаються з обмежень за X рядків вихідних таблиць.

Отже, $\pi_{X,R} : \Psi(R) \rightarrow \Psi(R \cap X)$, $\pi_{X,R}(\langle \psi, R \rangle) = \langle \psi', R \cap X \rangle$, де $\langle \psi, R \rangle \in \Psi(R)$.

Основа мультимножини ψ' визначається як $\Theta(\psi') = \{s \mid X \mid s \in \Theta(\psi)\}$.

Дублікати рядків, які з'явилися після виконання операції, не вилучаються. Кількість дублікатів кожного рядка визначається за формулою:

$$Occ(s', \psi') = \sum_{\substack{s \in \Theta(\psi), \\ s \upharpoonright X = s'}} Occ(s, \psi), \text{ де } s' \in \Theta(\psi').$$

Під з'єднанням таблиць схем R_1, R_2 розуміється бінарна параметрична операція \otimes_{R_1, R_2} , значеннями якої є таблиці схеми $R_1 \cup R_2$, що складаються з усяких об'єднань сумісних рядків вихідних таблиць.

Отже,

$$\otimes_{R_1, R_2} : \Psi(R_1) \times \Psi(R_2) \rightarrow \Psi(R_1 \cup R_2), \langle \psi_1, R_1 \rangle \otimes_{R_1, R_2} \langle \psi_2, R_2 \rangle = \langle \psi', R_1 \cup R_2 \rangle,$$

де $\langle \psi_1, R_1 \rangle \in \Psi(R_1)$, $\langle \psi_2, R_2 \rangle \in \Psi(R_2)$. Змістовно кажучи, кожний рядок з ψ_1 з'єднується з кожним рядком із ψ_2 , незалежно від того – дублікат це чи ні.

Основою мультимножини ψ' є множина рядків

$$\Theta(\psi') = \{s' \mid \exists s_1 \exists s_2 (s_1 \in \Theta(\psi_1) \wedge s_2 \in \Theta(\psi_2) \wedge s_1 \approx s_2 \wedge s' = s_1 \cup s_2)\}.$$

Кількість дублікатів знаходиться так: $Occ(s', \psi') = Occ(s' \upharpoonright R_1, \psi_1) \cdot Occ(s' \upharpoonright R_2, \psi_2)$, де $s' \in \Theta(\psi')$.

Введемо операцію перейменування. Під перейменуванням таблиць схеми R розуміється унарна параметрична операція $R\psi_{\xi,R}$, де $\xi : A \xrightarrow{\sim} A$ – ін'єктивне відображення на множині атрибутів, що здійснює перейменування атрибутів вихідних таблиць згідно з відображенням-параметром ξ . Перейменувати таблицю означає перейменувати атрибути її схеми, тобто перейменувати рядки таблиці. Перейменування рядків будемо здійснювати відповідно до [6].

Нехай $\eta : A \rightarrow A$ – функція перейменування атрибутів. Під перейменуванням рядків, що відповідає функції перейменування атрибутів η , розуміється відображення $Rs_{\eta} : S \rightarrow S'$, $Rs_{\eta}(s) = \{\langle \eta(A), s(A) \rangle \mid A \in \pi_1^2 s\}$, причому $\eta = \xi \cup \text{Id}_{A \setminus \text{dom } \xi}$.

Схема R називається ξ -допустимою, якщо $\xi[R] \cap (R \setminus \text{dom } \xi) = \emptyset$. Тут $\xi[R]$ – повний образ множини R відносно функції ξ . Множину таблиць схеми R , де схема R ξ -допустима позначимо $\Psi_{\xi}(R)$.

Під перейменуванням таблиць схеми R , що відповідає ін'єктивній частковій функції перейменування атрибутів $\xi: A \rightarrow A$, розуміється унарна параметрична операція $R\psi_{\xi,R}$ з областю означеності $\Psi_{\xi}(R)$, значення якої задаються таким чином: $R\psi_{\xi,R}(\langle\psi, R\rangle) = \langle Rs_{\eta}[\psi], \eta[R]\rangle$, де $\psi \in \Psi_{\xi}(R)$, $\eta = \xi \cup \text{id}_{A \setminus \text{dom}\xi}$, $Rs_{\eta}[\psi]$ – повний образ мультимножини ψ з основою $\Theta(\psi)$ відносно функції Rs_{η} .

Основою мультимножини $Rs_{\eta}[\psi]$ є повний образ множини $\Theta(\psi)$ відносно функції $Rs_{\eta,R}$. А кількість дублікатів рядка у результуючій таблиці задається рівністю

$$\text{Occ}(s', Rs_{\eta}[\psi]) = \text{Occ}(s, \psi),$$

де $s \in Rs_{\eta}^{-1}(s')$, $s' \in \Theta(\psi')$.

Введемо операцію активного доповнення. Для цього введемо декілька допоміжних понять.

Активним доменом атрибута $A \in R$ відносно таблиці $\langle\psi, R\rangle$ називається таблиця

$$D_{A,\psi} = \pi_{\{A\},R}(\langle\psi, R\rangle).$$

Насиченням таблиці $\langle\psi, R\rangle$ є таблиця

$$C(\langle\psi, R\rangle) = \pi_{\{A_1\},R}(\langle\psi, R\rangle) \otimes_{\{A_1\},\{A_2\}} \dots \otimes_{\{A_1, \dots, A_{n-1}\},\{A_n\}} \pi_{\{A_n\},R}(\langle\psi, R\rangle).$$

Під активним доповненням таблиці схеми R розуміється унарна параметрична операція \sim_R , яка таблиці зіставляє доповнення в її насиченні.

Отже, \sim_R :

$$\Psi(R) \rightarrow \Psi(R), \sim_R(\langle\psi, R\rangle) = C(\langle\psi, R\rangle) \setminus_{All}^R \langle\psi, R\rangle, \text{ де } \langle\psi, R\rangle \in \Psi(R).$$

Твердження. Мають місце наступні твердження:

- 1) будь-який вираз мультимножинної табличної алгебри можна замінити еквівалентним йому виразом, який використовує лише операції селекції, з'єднання, проєкції, об'єднання, різниці та перейменування;
- 2) будь-який вираз мультимножинної табличної алгебри, який містить лише скінченні таблиці можна замінити еквівалентним йому виразом, який використовує сталі таблиці з одним атрибутом і одним рядком, операції селекції, з'єднання, проєкції, об'єднання, різниці й перейменування.

Агрегатні операції

Широко використовуваними агрегатними операціями є *Sum*, *Avg*, *Min*, *Max*, *Count*. Так, операція *Sum* розраховує суму значень у відповідному стовпці заданої таблиці, при цьому значення *NULL* ігноруються. Операція *Avg* визначає середнє арифметичне значень у відповідному стовпці заданої таблиці, при цьому значення *NULL* ігноруються. Операції *Min* та *Max* знаходять найменше та найбільше значення у відповідному стовпці заданої таблиці, при цьому значення *NULL* так само ігноруються. Операція *Count* визначає кількість значень, відмінних від *NULL*, у відповідному стовпці заданої таблиці. Операція *Count(*)* визначає кількість рядків у заданій таблиці.

Нехай $\langle\psi, R\rangle \in \Psi(R)$, причому ψ – скінченна мультимножина і A – деякий атрибут схеми R , $A \in R$.

Позначимо через α_A – мультимножину, яка містить всі елементи стовпця з атрибутом A таблиці $\langle\psi, R\rangle$. Тоді $\alpha_A = D_{A,\psi}$, де $D_{A,\psi} = \pi_{A,R}(\langle\psi, R\rangle)$ – активний домен атрибута A відносно таблиці $\langle\psi, R\rangle$. Нехай $2_m^{D'} = \{\alpha \mid \Theta(\alpha) \in 2^{D'}\}$ – сім'я всіх мультимножин, основи яких є скінченними підмножинами множини D' ; тут $D' \subseteq D$ – підмножина універсального домену.

Нехай *Num* – числова підмножина універсального домену D , замкнена відносно додавання. Множина *Num* розширена включенням особливого елемента *NULL*, але при цьому операція додавання на випадок, коли хоча б один з аргументів є *NULL* не розширюється.

Задамо агрегатні операції *Sum*, *Avg*, *Min*, *Max*, *Count*. Їхніми аргументами є скінченні таблиці, а значеннями – одноатрибутні таблиці з одним рядком. Загальна схема: спочатку на скінченній мультимножині визначаються функції сумування, взяття найменшого та найбільшого значення, визначення середнього арифметичного і кількості елементів, а потім ці функції переносяться на таблиці.

Під операцією агрегування $Sum_{A,R}$ за атрибутом A (скінченних) таблиць схеми R , $A \in R$, розуміється унарна параметрична операція вигляду

$$Sum_{A,R} : \Psi(R) \rightarrow \Psi(\{A\}), \quad Sum_{A,R}(\langle \psi, R \rangle) = \left\langle \left\{ \left\{ A, Sum(\alpha_A) \right\} \right\}^1, \{A\} \right\rangle,$$

де $\langle \psi, R \rangle \in \Psi(R)$, а Sum – функція, що повертає суму значень стовпця з атрибутом A таблиці $\langle \psi, R \rangle$ (ці значення можуть повторюватися), які відрізняються від значення $NULL$, крім того, цей стовпець містить лише дані числового типу. Отже,

$$Sum(\alpha_A) = \begin{cases} NULL, & \text{якщо } \Theta(\alpha_A) = \emptyset; \\ NULL, & \text{якщо } \Theta(\alpha_A) = \{NULL\}; \\ \sum_{d \in \Theta(\alpha_A) \setminus \{NULL\}} d\alpha_A(d), & \text{якщо } \Theta(\alpha_A) \setminus \{NULL\} \neq \emptyset. \end{cases}$$

Верхній індекс 1 вказує на те, що рядок $\{A, Sum(\alpha_A)\}$ входить у результуючу таблицю лише один раз, тобто кількість дублікатів даного рядка дорівнює одиниці.

Таким чином,

$$Sum(\emptyset_m) = NULL, \quad Sum(\{NULL^n\}) = NULL, \quad Sum(\langle d_1^{n_1}, \dots, d_k^{n_k} \rangle) = \sum_{i=1}^k d_i n_i,$$

в припущенні, що всі елементи d_i відрізняються від елемента $NULL$.

Для випадку порожньої таблиці $\langle \psi_{\emptyset}, R \rangle$ операція агрегування $Sum_{A,R}$ застосовується так:

$$Sum_{A,R}(\langle \psi_{\emptyset}, R \rangle) = \left\langle \left\{ \left\{ A, NULL \right\} \right\}^1, \{A\} \right\rangle.$$

Приклад 1. Серед студентів фізико-математичного факультету було проведено анонімне опитування на визначення рівня самооцінки. Проаналізовано відповіді 6 учасників вибраних довільним чином з кожного факультету. Дані відображені в таблиці $\langle \Phi M, R \rangle$, де $R = \{A, B, C, D\}$ (див. таблицю). Атрибути схеми R – це запитання, дані в таблицях – це відповіді студентів.

Запитання (атрибути):

A – чи турбуєтесь Ви за свій психічний стан?

B – чи турбуєтесь Ви про своє майбутнє?

C – чи боїтесь Ви виступати з промовою перед незнайомими людьми?

D – чи часто Ви робите помилки?

Варіанти відповіді:

1 – часто;

2 – іноді;

3 – рідко;

4 – ніколи.

Таблиця

A	B	C	D
4	2	4	3
4	2	4	3
4	2	4	3
4	3	3	2
4	3	3	2
3	3	3	1

Позначимо рядки:

$$s_1 = \{ \langle A, 4 \rangle, \langle B, 2 \rangle, \langle C, 4 \rangle, \langle D, 3 \rangle \}, \quad s_2 = \{ \langle A, 4 \rangle, \langle B, 3 \rangle, \langle C, 3 \rangle, \langle D, 2 \rangle \},$$

$$s_3 = \{ \langle A, 3 \rangle, \langle B, 3 \rangle, \langle C, 3 \rangle, \langle D, 1 \rangle \},$$

тоді $\Phi M = \{s_1^3, s_2^2, s_3^1\}$.

Тоді в результаті застосування операції агрегування $Sum_{A,R}$ до таблиці $\langle \Phi M, R \rangle$ за атрибутами A, B, C, D отримуємо таблиці:

$$Sum_{A,R} = \langle \{\{A, 23\}\}, \{A\} \rangle, \quad Sum_{B,R} = \langle \{\{B, 15\}\}, \{B\} \rangle,$$

$$Sum_{C,R} = \langle \{\{C, 21\}\}, \{C\} \rangle, \quad Sum_{D,R} = \langle \{\{D, 14\}\}, \{D\} \rangle.$$

Нехай \leq – лінійний порядок на універсальному домені D . Під операцією агрегування $Min_{A,R}$ за атрибутом A (скінченних) таблиць схеми R , $A \in R$, розуміється унарна параметрична операція вигляду:

$$Min_{A,R} : \Psi(R) \rightarrow \Psi(\{A\}), \quad \text{де } Min_{A,R}(\langle \psi, R \rangle) = \langle \{\{A, Min(\alpha_A)\}\}^{\uparrow}, \{A\} \rangle,$$

причому $\langle \psi, R \rangle \in \Psi(R)$, а Min – функція, що повертає найменше значення серед значень стовпця з атрибутом A таблиці $\langle \psi, R \rangle$, відмінних від значення $NULL$, тобто $Min : 2_m^D \rightarrow D$,

$$Min(\alpha_A) = \begin{cases} NULL, & \text{якщо } \Theta(\alpha_A) = \emptyset; \\ NULL, & \text{якщо } \Theta(\alpha_A) = \{NULL\}; \\ \min\{d \mid d \in \Theta(\alpha_A) \setminus \{NULL\}\}, & \text{якщо } \Theta(\alpha_A) \setminus \{NULL\} \neq \emptyset. \end{cases}$$

Таким чином, $Min(\emptyset_m) = NULL$, $Min(\{NULL^n\}) = NULL$, $Min(\{d_1^{n_1}, \dots, d_k^{n_k}\}) = \min\{d_1, \dots, d_k\}$, в припущенні, що всі елементи d_i , $i = \overline{1, k}$, відрізняються від елемента $NULL$.

Для випадку порожньої таблиці $\langle \psi_\emptyset, R \rangle$ операція агрегування $Min_{A,R}$ застосовується так:

$$Min_{A,R}(\langle \psi_\emptyset, R \rangle) = \langle \{\{A, NULL\}\}^{\uparrow}, \{A\} \rangle.$$

Приклад 2. Розглянемо таблицю із приклада 1. Застосуємо операцію агрегування $Min_{A,R}$ до таблиці $\langle \Phi M, R \rangle$ за атрибутами A, B, C, D отримуємо таблиці:

$$Min_{A,R} = \langle \{\{A, 3\}\}, \{A\} \rangle, \quad Min_{B,R} = \langle \{\{B, 2\}\}, \{B\} \rangle,$$

$$Min_{C,R} = \langle \{\{C, 3\}\}, \{C\} \rangle, \quad Min_{D,R} = \langle \{\{D, 1\}\}, \{D\} \rangle.$$

Під операцією агрегування $Max_{A,R}$ за атрибутом A (скінченних) таблиць схеми R , $A \in R$, розуміється унарна параметрична операція вигляду:

$$Max_{A,R} : \Psi(R) \rightarrow \Psi(\{A\}), \quad \text{де } Max_{A,R}(\langle \psi, R \rangle) = \langle \{\{A, Max(\alpha_A)\}\}^{\uparrow}, \{A\} \rangle,$$

де $\langle \psi, R \rangle \in \Psi(R)$, а Max – функція, що повертає найбільше значення серед значень стовпця з атрибутом A таблиці $\langle \psi, R \rangle$, відмінних від $NULL$, тобто $Max : 2_m^D \rightarrow D$,

$$Max(\alpha_A) = \begin{cases} NULL, & \text{якщо } \Theta(\alpha_A) = \emptyset; \\ NULL, & \text{якщо } \Theta(\alpha_A) = \{NULL\}; \\ \max\{d \mid d \in \Theta(\alpha_A) \setminus \{NULL\}\}, & \text{якщо } \Theta(\alpha_A) \setminus \{NULL\} \neq \emptyset. \end{cases}$$

Таким чином, $Max(\emptyset_m) = NULL$, $Max(\{NULL^n\}) = NULL$, $Max(\{d_1^{n_1}, \dots, d_k^{n_k}\}) = \max\{d_1, \dots, d_k\}$, в припущенні, що всі елементи d_i , $i = \overline{1, k}$, відрізняються від значення $NULL$.

Для випадку порожньої таблиці $\langle \psi_\emptyset, R \rangle$ операція агрегування $Max_{A,R}$ застосовується так:

$$Max_{A,R}(\langle \psi_\emptyset, R \rangle) = \langle \{\{A, NULL\}\}^{\uparrow}, \{A\} \rangle.$$

Приклад 3. Застосуємо операцію агрегування $Max_{A,R}$ до таблиці $\langle \Phi M, R \rangle$ за атрибутами A, B, C, D отримуємо таблиці:

$$Max_{A,R} = \langle \{ \{A, 4\} \}, \{A\} \rangle, Max_{B,R} = \langle \{ \{B, 3\} \}, \{B\} \rangle,$$

$$Max_{C,R} = \langle \{ \{C, 4\} \}, \{C\} \rangle, Max_{D,R} = \langle \{ \{D, 3\} \}, \{D\} \rangle.$$

Значимо, що функції Min та Max визначають найменший або найбільший елемент основи множини α_A серед елементів основи, відмінних від $NULL$, тому порівнянність особливого елемента $NULL$ з рештою елементів універсального домену в даному випадку неістотна.

Під операцією агрегування $Count_{A,R}$ за атрибутом A (скінченних) таблиць схеми R , $A \in R$, розуміється унарна параметрична операція вигляду:

$$Count_{A,R} : \Psi(R) \rightarrow \Psi(\{A\}), Count_{A,R}(\langle \psi, R \rangle) = \langle \{ \{A, Count(\alpha_A)\} \}^1, \{A\} \rangle,$$

де $\langle \psi, R \rangle \in \Psi(R)$, а $Count$ – функція, що повертає кількість значень, відмінних від $NULL$, з урахуванням дублікатів, у стовпці з атрибутом A таблиці $\langle \psi, R \rangle$, тобто:

$$Count : 2_m^D \rightarrow Z_+, Count(\alpha_A) = \sum_{d \in \Theta(\alpha_A) \setminus \{NULL\}} \alpha_A(d);$$

покладається за означенням, що сума порожньої множини доданків дорівнює нулю.

Таким чином, $Count(\emptyset_m) = 0$, $Count(\{NULL^n\}) = 0$, $Count(\{d_1^{n_1}, \dots, d_k^{n_k}\}) = n_1 + \dots + n_k$, в припущенні, що всі елементи d_i , $i = \overline{1, k}$, відрізняються від елемента $NULL$.

Для випадку порожньої таблиці $\langle \psi_{\emptyset}, R \rangle$ операція агрегування $Count_{A,R}$ застосовується так:

$$Count_{A,R}(\langle \psi_{\emptyset}, R \rangle) = \langle \{ \{A, 0\} \}^1, \{A\} \rangle.$$

Приклад 4. Застосуємо операцію агрегування $Count_{A,R}$ до таблиці $\langle \Phi M, R \rangle$ за атрибутами A, B, C, D отримаємо таблиці:

$$Count_{A,R} = \langle \{ \{A, 6\} \}, \{A\} \rangle, Count_{B,R} = \langle \{ \{B, 6\} \}, \{B\} \rangle,$$

$$Count_{C,R} = \langle \{ \{C, 6\} \}, \{C\} \rangle, Count_{D,R} = \langle \{ \{D, 6\} \}, \{D\} \rangle.$$

Припустимо, що числова підмножина Num універсального домену замкнена відносно (часткової операції) ділення $/ : Num \times Num \rightarrow Num$. Довизначимо операцію ділення так, що коли перший аргумент дорівнює $NULL$, то функція приймає значення $NULL$.

Під операцією агрегування $Avg_{A,R}$ за атрибутом A (скінченних) таблиць схеми R , $A \in R$, розуміється унарна параметрична операція вигляду:

$$Avg_{A,R} : \Psi(R) \rightarrow \Psi(\{A\}), Avg_{A,R}(\langle \psi, R \rangle) = \langle \{ \{A, Avg(\alpha_A)\} \}^1, \{A\} \rangle,$$

де $\langle \psi, R \rangle \in \Psi(R)$, а Avg – функція, що повертає середнє арифметичне значення елементів стовпця з атрибутом A таблиці $\langle \psi, R \rangle$, які відрізняються від значення $NULL$, з урахуванням дублікатів, тобто:

$$Avg : 2_m^{Num} \rightarrow Num, Avg(\alpha_A) = Sum(\alpha_A) / Count(\alpha_A).$$

Таким чином, з означення випливають рівності

$$Avg(\emptyset_m) = Sum(\emptyset_m) / Count(\emptyset_m) = NULL / 0 = NULL,$$

$$Avg(\{NULL^n\}) = Sum(\{NULL^n\}) / Count(\{NULL^n\}) = NULL / 0 = NULL,$$

$$Avg(\{d_1^{n_1}, \dots, d_k^{n_k}\}) = Sum(\{d_1^{n_1}, \dots, d_k^{n_k}\}) / Count(\{d_1^{n_1}, \dots, d_k^{n_k}\}) = \frac{1}{(n_1 + \dots + n_k)} \sum_{i=1}^k d_i n_i,$$

в припущенні, що всі елементи d_i відрізняються від елемента $NULL$.

Для випадку порожньої таблиці $\langle \psi_{\emptyset}, R \rangle$ операція агрегування $Avg_{A,R}$ застосовується так:

$$Avg_{A,R}(\langle \psi_{\emptyset}, R \rangle) = \left\langle \left\{ \left\{ \left\{ A, NULL \right\} \right\}^{\parallel} \right\}, \{A\} \right\rangle.$$

Приклад 5. Застосуємо операцію агрегування $Avg_{A,R}$ до таблиці $\langle \Phi M, R \rangle$ за атрибутами A, B, C, D отримаємо таблиці

$$Avg_{A,R} = \left\langle \left\{ \left\{ A, \frac{23}{6} \right\} \right\}, \{A\} \right\rangle, \quad Avg_{B,R} = \left\langle \left\{ \left\{ B, \frac{15}{6} \right\} \right\}, \{B\} \right\rangle,$$

$$Avg_{C,R} = \left\langle \left\{ \left\{ C, \frac{21}{6} \right\} \right\}, \{C\} \right\rangle, \quad Avg_{D,R} = \left\langle \left\{ \left\{ D, \frac{14}{6} \right\} \right\}, \{D\} \right\rangle.$$

Під операцією агрегування $Count_{A,R}(\ast)$ (скінченних) таблиць схеми R розуміється унарна параметрична операція вигляду:

$$Count_{A,R}(\ast) : \Psi(R) \rightarrow \Psi(\{A\}), \quad Count_{A,R}(\ast)(\langle \psi, R \rangle) = \left\langle \left\{ \left\{ \left\{ A, \|\psi\| \right\} \right\}^{\parallel} \right\}, \{A\} \right\rangle,$$

при цьому $\langle \psi, R \rangle \in \Psi(R)$, а $\|\psi\|$ – це ранг мультимножини ψ .

Для випадку порожньої таблиці $\langle \psi_{\emptyset}, R \rangle$ операція агрегування $Count_{A,R}(\ast)$ застосовується так:

$$Count_{A,R}(\ast)(\langle \psi_{\emptyset}, R \rangle) = \left\langle \left\{ \left\{ \left\{ A, \|\emptyset_m\| \right\} \right\}^{\parallel} \right\}, \{A\} \right\rangle = \left\langle \left\{ \left\{ \left\{ A, 0 \right\} \right\}^{\parallel} \right\}, \{A\} \right\rangle.$$

Приклад 6. Застосуємо операцію агрегування $Count_{A,R}(\ast)$ до таблиці $\langle \Phi M, R \rangle$ за атрибутами A, B, C, D отримаємо таблиці:

$$Count_{A,R}(\ast) = \left\langle \left\{ \left\{ \left\{ A, 6 \right\} \right\}, \{A\} \right\rangle, \quad Count_{B,R}(\ast) = \left\langle \left\{ \left\{ \left\{ B, 6 \right\} \right\}, \{B\} \right\rangle,$$

$$Count_{C,R}(\ast) = \left\langle \left\{ \left\{ \left\{ C, 6 \right\} \right\}, \{C\} \right\rangle, \quad Count_{D,R}(\ast) = \left\langle \left\{ \left\{ \left\{ D, 6 \right\} \right\}, \{D\} \right\rangle.$$

Висновки

Визначено формальну математичну семантику агрегатних операцій над таблицями як мультимножинами рядків однієї схеми, яка проілюстрована прикладами їх використання. Для заданих агрегатних операцій параметром виступає лише один атрибут, проте отримані результати можна розширити на випадки, коли параметром є деяка функція над рядком. Результати роботи можуть бути використані в теорії узагальнених табличних алгебр.

1. Codd E.F. A Relational Model of Data for Large Shared Data Banks / E.F. Codd // Comm. of ACM. – 1970. – 13, N 6. – P. 377–387.
2. Grefen Paul W.P.J., Rolf A. De By A Multi-Set Extended Relational Algebra. A Formal Approach to a Practical Issue // 10th International Conference on Data Engineering, ICDE, February 14-18, 1994, Houston, TX, USA. – 1994. – P. 80–88.
3. Lamperti G., Melchiori M., Zanella M. On Multisets in Database Systems // Multiset Processing: Mathematical, Computer Science, and Molecular Computing Points of View, number 2235 in Lecture Notes in Computing Since. – Berlin: Springer-Verlag. – 2001. – P. 147–215.
4. Garcia-Molina H. Database Systems: The Complete Book: [2nd Edition] / H. Garcia-Molina, J.D. Ullman, J. Widom. – Prentice Hall, 2008. – 1119 p.
5. Silbeschatz A., Korth H., Sudarshan S. Database System Concepts – McGraw-Hill, 2011. – 1376 p.
6. Редько В.Н., Брона Ю.Й., Буй Д.Б., Поляков С.А. Реляційні бази даних: табличні алгебри та SQL-подібні мови. – Київ: Видавничий дім "Академперіодика", 2001. – 198 с.