

УДК 004.655

Мультимножественная табличная алгебра: дополнительные операции

Д. Б. Буй*, И. Н. Глушко**

* Киевский национальный университет имени Тараса Шевченко, Киев, Украина

** Нежинский государственный университет имени Николая Гоголя, Нежин, Украина

Сигнатура мультимножественной табличной алгебры пополнена новыми операциями: операциями внутренних и внешних соединений, операцией полусоединения, агрегатными операциями. Задана формальная математическая семантика указанных операций. Для задания внешних соединений введен особый элемент универсального домена NULL.

Ключевые слова: реляционные базы данных, мультимножественная табличная алгебра, расширенная мультимножественная табличная алгебра.

Сигнатура мультимножественной табличной алгебры пополнена новыми операциями: операциями внутренних и внешних соединений, операцией полусоединения, агрегатными операциями. Задана формальная математическая семантика указанных операций. Для задания внешних соединений введен особый элемент универсального домена NULL.

Ключеві слова: реляційні бази даних, мультимножественная табличная алгебра, розширена мультимножественная табличная алгебра.

The signature of multiset table algebra is filled up with new operations such as inner and outer joins, semijoin and aggregate operations. A formal mathematical semantics of these operations is defined. The special element NULL is introduced into the universal domain to define of outer join.

Key words: relation databases, multiset table algebra, extending multiset table algebra.

1. Общая постановка задачи и её актуальность

Реляционная модель данных в настоящее время широко используется как в научных исследованиях в базах данных, так и на практике. Данная модель основана на множествах кортежей, то есть не позволяет дубликаты кортежей в отношении [1]. Однако многие языки, ориентированные на работу с базами данных, требуют реляционную модель данных с мультимножественной семантикой (multiset semantics). Это предполагает понимание таблиц как мультимножеств, т.е. совокупностей с дубликатами. Вопросу использования мультимножеств в базах данных уделяли внимание G. Lamperti, M. Melchiori, M. Zanella [2], Г. Гарсиа-Молина, Дж. Ульман, Дж. Уидом [3], А. Silbeschatz, Н. Korth, S. Sudarshan [4], а также отечественные ученые Д.Б. Буй, С.А. Поляков, Ю.Й. Брона, В.Н. Редько [5]. Обзор литературы об использовании мультимножеств в базах данных проведен в статье [6], которая насчитывает 9 источников по данной теме.

Вместе с тем, этот вопрос требует уточнения и расширения, поскольку ни в одной из указанных работ не уделяется достаточное внимание операциям внутренних и внешних соединений, операции полусоединения, внешним мультимножественным операциям, а также агрегатным операциям над таблицами мультимножественной табличной алгебры.

2. Основные понятия теории мультимножеств

Приведем основные понятия мультимножеств в терминах работ [5, 7]. Зафиксируем множество U . Под мультимножеством α с основой U понимаем отображение вида $\alpha : U \rightarrow \{1, 2, \dots\}$. Пусть D – универсум элементов основ мультимножеств, тогда булеан $P(D)$ – универсум основ мультимножеств.

Под характеристической функцией мультимножества α понимаем функцию вида $\chi_\alpha : D \rightarrow \{0, 1, 2, \dots\}$, значение которой задается кусочной схемой :

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

для всех $d \in D$, где $\text{dom } \alpha$ – область определения мультимножества α , т.е. его основа.

Мультимножество называется пустым и обозначается как \emptyset_m , если его основа – пустое множество.

Мультимножества, областью значений которых является пустое множество или одноэлементное множество вида $\{1\}$ называются 1-мультимножествами. Эти мультимножества есть аналогами обычных множеств.

Договоримся мультимножество α с основой $\{d_1, \dots, d_k\}$ записывать как $\{d_1^{n_1}, \dots, d_k^{n_k}\}$, где n_i – количество дубликатов (экземпляров) элемента d_i в мультимножестве α , т.е. $n_i = \alpha(d_i)$, $i = 1, \dots, k$.

Под рангом конечного мультимножества α понимаем количество дубликатов элементов его основы $\|\alpha\| = \sum_{d \in \text{dom } \alpha} \alpha(d)$; при этом $\|\emptyset_m\| = 0$.

Скажем, что мультимножество β включается в мультимножество α ($\beta \preceq \alpha$), если: $\beta \preceq \alpha \Leftrightarrow U_\beta \subseteq U_\alpha \ \& \ \forall d (d \in U_\beta \Rightarrow \beta(d) \leq \alpha(d))$. Здесь U_α и U_β основы мультимножеств α и β соответственно.

Если $\beta \preceq \alpha$, то мультимножество β называется подмультимножеством мультимножества α , а мультимножество α – надмультимножеством мультимножества β .

В работе [5] операции над мультимножествами определены в терминах характеристических функций. Авторы определяют операции объединения \cup_1 , пересечения \cap_1 , разности \setminus_1 мультимножеств, которые строят 1-мультимножества, основы которых получаются соответственно теоретико-множественными объединением, пересечением и разницей основ мультимножеств-аргументов. Кроме того, вводятся операции объединения \cup_{All} , пересечения \cap_{All} , разности \setminus_{All} мультимножеств, которые строят мультимножества общего вида. Также задано операцию декартового соединения мультимножеств \otimes и операцию $Dist(\alpha)$, которая строит 1-мультимножество, основа которого совпадает с основой исходного мультимножества. Наконец, в этой работе вводится аналог полного образа (множества относительно функции) для мультимножеств.

3. Мультимножественная табличная алгебра

Рассмотрим два множества: A – множество атрибутов (имен) и D – универсальный домен (множество денотатов). Произвольное (конечное) множество атрибутов $R \subseteq A$ назовём схемой. Под строкой схемы R понимаем именованное множество на паре A и D [5], проекция которого по первой компоненте совпадает с R , т. е. строка схемы R – это функция вида $s: R \rightarrow D$. Множество всех строк схемы R обозначим $S(R)$, а множество всех строк – S .

Под таблицей схемы R понимаем пару $\langle \psi, R \rangle$, где первая компонента ψ – это произвольное мультимножество, основой которого $\Theta(\psi)$ является произвольное множество, в частности, бесконечное, строк схемы R , а вторая компонента R – схема (таблицы).

Под мультимножественной табличной алгеброй понимаем алгебру $\langle \Psi, \Omega_{P, \Xi} \rangle$, где $\Psi = \bigcup_{R \subseteq A} \Psi(R)$ – множество всех таблиц, $\Psi(R)$ – множество всех таблиц схемы R ,

$$\Omega_{P, \Xi} = \left\{ \bigcup_{All}^R, \bigcap_{All}^R, \setminus_{All}^R, \sigma_{p, R}, \pi_{X, R}, \otimes_{R_1, R_2}, Rt_{\xi, R}, \sim_R \right\}_{\substack{p \in P, \xi \in \Xi \\ X, R, R_1, R_2 \subseteq A}} \text{ – сигнатура; } P,$$

Ξ – множества параметров. Операции мультимножественной табличной алгебры задано в [8].

Пополним сигнатуру мультимножественной табличной алгебры новыми операциями: операциям внутренних и внешних соединений, операцией полусоединения, агрегатными операциями.

4. Операции внутреннего соединения

Под декартовым соединением C_j (Cartesian Join) таблиц схем R_1 и R_2 , причем $R_1 \cap R_2 = \emptyset$, понимаем бинарную параметрическую операцию вида

$$C_j : \Psi(R_1) \times \Psi(R_2) \rightarrow \Psi(R_1 \cup R_2), \quad \langle \psi_1, R_1 \rangle_{R_1, R_2} C_j \langle \psi_2, R_2 \rangle = \langle \psi', R_1 \cup R_2 \rangle, \quad \text{где}$$

$\langle \psi_1, R_1 \rangle \in \Psi(R_1)$, $\langle \psi_2, R_2 \rangle \in \Psi(R_2)$. Основой мультимножества ψ' является множество строк $\Theta(\psi') = \{s \mid \exists s_1 \exists s_2 (s_1 \in \Theta(\psi_1) \wedge s_2 \in \Theta(\psi_2) \wedge s = s_1 \cup s_2)\}$. Количество дубликатов определяется так: $Occ(s, \psi') = Occ(s_1, \psi_1) \cdot Occ(s_2, \psi_2)$, где $s \in \Theta(\psi')$ и $s = s_1 \cup s_2$ (очевидно, что это разложение строки s единственно).

Ниже используется бинарное отношение совместимости строк, которое вводится так: $s_1 \approx s_2 \stackrel{def}{\Leftrightarrow} s_1 \upharpoonright R' = s_2 \upharpoonright R'$, де $R' = R_1 \cap R_2$, а R_1, R_2 – схемы строк s_1, s_2 соответственно [5]. Основное свойство этого отношения состоит в следующем: $s_1 \cup s_2 \in S(R_1 \cup R_2) \Leftrightarrow s_1 \approx s_2$ [5].

Под внутренним естественным соединением (Inner Natural Join) таблиц схем R_1 и R_2 понимаем бинарную параметрическую операцию \otimes_{R_1, R_2} , значениями

которой являются таблицы схемы $R_1 \cup R_2$, которые, говоря содержательно, содержат все объединения совместных строк исходных таблиц. Таким образом, $\otimes_{R_1, R_2} : \Psi(R_1) \times \Psi(R_2) \rightarrow \Psi(R_1 \cup R_2)$, $\langle \psi_1, R_1 \rangle \otimes_{R_1, R_2} \langle \psi_2, R_2 \rangle = \langle \psi', R_1 \cup R_2 \rangle$, где $\langle \psi_1, R_1 \rangle \in \Psi(R_1)$, $\langle \psi_2, R_2 \rangle \in \Psi(R_2)$. Основой мультимножества ψ' является множество строк вида

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

Количество дубликатов определяется, как и ранее, так: $Occ(s, \psi') = Occ(s_1, \psi_1) \cdot Occ(s_2, \psi_2)$, где $s \in \Theta(\psi')$ и $s = s_1 \cup s_2$ (как и ранее, приведенное разложение единственно).

Под внутренним соединением по атрибутам A_1, \dots, A_n (Inner Join using A_1, \dots, A_n), причем все A_1, \dots, A_n попарно различны, $n \geq 1$, таблиц схем R_1 и R_2 , где $R_1 \cap R_2 \supseteq \{A_1, \dots, A_n\}$ (подразумевается, что общие атрибуты, отличающиеся от атрибутов A_1, \dots, A_n , перед соединением переименовываются) понимаем бинарную параметрическую операцию вида

$$\otimes_{A_1, \dots, A_n, R_1, R_2} : \Psi(R_1) \times \Psi(R_2) \rightarrow \Psi(R_1 \cup R_2),$$

причем $\langle \psi_1, R_1 \rangle \otimes_{A_1, \dots, A_n, R_1, R_2} \langle \psi_2, R_2 \rangle = \langle \psi', R_1 \cup R_2 \rangle$, где $\langle \psi_1, R_1 \rangle \in \Psi(R_1)$,

$\langle \psi_2, R_2 \rangle \in \Psi(R_2)$. Основой мультимножества ψ' является множество строк

$$\Theta(\psi') = \left\{ s \mid \exists s_1 \exists s_2 \left(s_1 \in \Theta(\psi_1) \wedge s_2 \in \Theta(\psi_2) \wedge \bigwedge_{i=1}^n s_1(A_i) = s_2(A_i) \wedge s = s_1 \cup s_2 \right) \right\}.$$

Количество дубликатов определяется, как и ранее, так: $Occ(s, \psi') = Occ(s_1, \psi_1) \cdot Occ(s_2, \psi_2)$, где $s \in \Theta(\psi')$ и $s = s_1 \cup s_2$ (как и ранее, приведенное разложение единственно).

Отметим, если таблицы-аргументы имеют еще и другие общие атрибуты, отличающиеся от атрибутов A_1, \dots, A_n , то перед соединением их нужно переименовать.

Пусть $p : S \times S \rightarrow \{true, false\}$ – вообще говоря частичный бинарный предикат на множестве всех строк S , такой, что выполняется импликация $\forall s_1 \forall s_2 ((s_1, s_2) \in \text{dom } p \wedge p(s_1, s_2) = true \Rightarrow s_1 \approx s_2)$.

Под внутренним соединением по предикату p (Inner Join on p) таблиц схем R_1 и R_2 понимаем частичную бинарную параметрическую операцию вида

$$\otimes_{p, R_1, R_2} : \Psi(R_1) \times \Psi(R_2) \rightarrow \Psi(R_1 \cup R_2), \quad \text{dom } \otimes_{p, R_1, R_2} = \{ \langle \psi_1, R_1 \rangle, \langle \psi_2, R_2 \rangle \mid \Theta(\psi_1) \times \Theta(\psi_2) \subseteq \text{dom } p \},$$

$$\langle \psi_1, R_1 \rangle \otimes_{p, R_1, R_2} \langle \psi_2, R_2 \rangle = \langle \psi', R_1 \cup R_2 \rangle. \quad \text{Основой}$$

мультимножества ψ' является множество строк $\Theta(\psi') = \{s \mid \exists s_1 \exists s_2 (s_1 \in \Theta(\psi_1) \wedge s_2 \in \Theta(\psi_2) \wedge p(s_1, s_2) \simeq true \wedge s = s_1 \cup s_2)\}$.

Количество дубликатов определяется, как и ранее, так: $Occ(s, \psi') = Occ(s_1, \psi_1) \cdot Occ(s_2, \psi_2)$, где $s \in \Theta(\psi')$ и $s = s_1 \cup s_2$. Выше предполагалось, что пара таблиц $\langle \psi_1, R_1 \rangle, \langle \psi_2, R_2 \rangle$ принадлежит указанной области определенности.

Отметим следующий очевидный факт. Операция естественного соединения \otimes_{R_1, R_2} является расширением произвольной другой операции соединения в

следующем смысле: $\langle \psi_1, R_1 \rangle \underset{R_1, R_2}{Cj} \langle \psi_2, R_2 \rangle = \langle \psi_1, R_1 \rangle \otimes_{R_1, R_2} \langle \psi_2, R_2 \rangle$,

$$\langle \psi_1, R_1 \rangle \otimes_{A_1, \dots, A_n, R_1, R_2} \langle \psi_2, R_2 \rangle = \langle \psi_1, R_1 \rangle \otimes_{R_1, R_2} \langle \psi_2, R_2 \rangle,$$

$$\left(\langle \psi_1, R_1 \rangle \otimes_{p, R_1, R_2} \langle \psi_2, R_2 \rangle \right)_1 \preceq \left(\langle \psi_1, R_1 \rangle \otimes_{R_1, R_2} \langle \psi_2, R_2 \rangle \right)_1^1,$$

при условии, что значения операций в левых частях этих двух равенств и включения определены.

Под операцией полусоединения (Semijoin) двух таблиц схем R_1 и R_2 понимаем бинарную параметрическую операцию \square_{R_1, R_2} , значением которой является таблица схемы R_1 , которая содержит те строки первой таблицы, которые входят в (естественное) соединение таблиц-аргументов.

Следовательно, $\square_{R_1, R_2} : \Psi(R_1) \times \Psi(R_2) \rightarrow \Psi(R_1)$, $\langle \psi_1, R_1 \rangle \square_{R_1, R_2} \langle \psi_2, R_2 \rangle = \langle \psi', R_1 \rangle$, где $\langle \psi_1, R_1 \rangle \in \Psi(R_1)$, $\langle \psi_2, R_2 \rangle \in \Psi(R_2)$. Основой мультимножества ψ' является множество строк $\Theta(\psi') = \{s_1 \mid s_1 \in \Theta(\psi_1) \wedge \exists s_2 (s_2 \in \Theta(\psi_2) \wedge s_1 \approx s_2)\}$. Количество дубликатов определяется так: $Occ(s, \psi') = Occ(s, \psi_1)$, где $s \in \Theta(\psi')$.

5. Операции внешнего соединения

При применении операций внутреннего соединения возможна потеря информации, поскольку строки одной таблицы, которые не соединяются со строками другой таблицы, не будут включены в результирующую таблицу. В тех случаях, когда необходимо учесть строки таблиц-аргументов, которые не попали в результат исходного внутреннего соединения, используют операции внешнего соединения.

Для обозначения отсутствующих значений атрибутов строк результирующей таблицы используем особый элемент универсального домена $NULL$. Обозначим как $s_{R, NULL}$ константную строку схемы R вида $s_{R, NULL} : R \rightarrow \{NULL\}$, присваивающую всем атрибутам своей схемы значения $NULL$.

¹Запись $(\langle \psi, R \rangle)_1$ обозначает первую компоненту пары $\langle \psi, R \rangle$, то есть мультимножество ψ .

Используем логическую схему определения операций внешнего соединения из [5], следуя которой четыре операции внешнего соединения вводятся как операции, подчиненные одной операции внутреннего соединения.

Пусть $\varphi: \Psi(R_1) \times \Psi(R_2) \rightarrow \Psi(R_1 \cup R_2)$ – некоторая частичная бинарная операция на множестве таблиц, причем выполняется включение $(\varphi(\langle \psi_1, R_1 \rangle, \langle \psi_2, R_2 \rangle))_1 \subseteq \left(\langle \psi_1, R_1 \rangle \otimes_{R_1, R_2} \langle \psi_2, R_2 \rangle \right)_1$ для всех $\langle \langle \psi_1, R_1 \rangle, \langle \psi_2, R_2 \rangle \rangle \in \text{dom } \varphi$.

Отметим, что операции внутреннего соединения C_j , \otimes_{R_1, R_2} , $\otimes_{A_1, \dots, A_n, R_1, R_2}$,

\otimes_{p, R_1, R_2} именно такие.

Зафиксируем таблицы $\langle \psi_1, R_1 \rangle, \langle \psi_2, R_2 \rangle$ из области определенности операции φ . Тогда таблица $\langle \psi_1, R_1 \rangle$ предполагает следующее представление:

$\langle \psi_1, R_1 \rangle = \left\langle \psi_1 \cap_{\varphi} \psi_2, R_1 \right\rangle \cup_{All}^{R_1} \left\langle \psi_1 - \psi_2, R_1 \right\rangle$, где $\left\langle \psi_1 \cap_{\varphi} \psi_2, R_1 \right\rangle = \langle \psi', R_1 \rangle$, основой мультимножества ψ' является множество строк

$$\Theta(\psi') = \{s_1 \mid s_1 \in \Theta(\psi_1) \wedge \exists s_2 (s_2 \in \Theta(\psi_2) \wedge s_1 \cup s_2 \in \Theta(\langle \langle \psi_1, R_1 \rangle, \langle \psi_2, R_2 \rangle \rangle))\},$$

а количество дубликатов $Occ(s_1, \psi') = Occ(s_1, \psi_1)$, $s_1 \in \Theta(\psi')$ и

$\left\langle \psi_1 - \psi_2, R_1 \right\rangle = \langle \psi'', R_1 \rangle$, основой мультимножества ψ'' является множество строк

$$\Theta(\psi'') = \{s_1 \mid s_1 \in \Theta(\psi_1) \wedge \forall s_2 (s_2 \in \Theta(\psi_2) \Rightarrow s_1 \cup s_2 \notin \Theta(\langle \langle \psi_1, R_1 \rangle, \langle \psi_2, R_2 \rangle \rangle))\},$$

а количество дубликатов определяется так: $Occ(s_1, \psi'') = Occ(s_1, \psi_1)$, где $s_1 \in \Theta(\psi'')$.

Говоря содержательно, строки таблицы $\left\langle \psi_1 \cap_{\varphi} \psi_2, R_1 \right\rangle$ используются при формировании результата (внутреннего) соединения, а строки таблицы $\left\langle \psi_1 - \psi_2, R_1 \right\rangle$ – не используются. Аналогичное представление таблицы $\langle \psi_2, R_2 \rangle$ получим, поменяв роли таблиц $\langle \psi_1, R_1 \rangle$ и $\langle \psi_2, R_2 \rangle$ в представлении таблицы $\langle \psi_1, R_1 \rangle$.

Отметим, что если операция φ совпадает с операцией \otimes_{R_1, R_2} , то

$\left\langle \psi_1 \cap_{\varphi} \psi_2, R_1 \right\rangle = \langle \psi_1, R_1 \rangle \square_{R_1, R_2} \langle \psi_2, R_2 \rangle$, т.е. таблица в левой части последнего

равенства получается в результате применения операции полусоединения к исходным таблицам.

Ниже для упрощения записи будем считать, что операции соединения имеют больший приоритет, чем операции объединения.

Определим четыре операции внешнего соединения (Outer Join), индуцированные одной операцией внутреннего соединения φ . Для этого рассмотрим следующие естественные соединения:

$$\left\langle \psi_1 - \psi_2, R_1 \right\rangle_{\varphi, R_1, R_2 \setminus R_1} \otimes \left\langle \{s_{R_2 \setminus R_1, NULL}^1\}, R_2 \setminus R_1 \right\rangle = \langle \psi', R_1 \cup R_2 \rangle,$$

где основой мультимножества ψ' является множество строк $\Theta(\psi') = \{s_1 \cup s_{R_2 \setminus R_1, NULL} \mid s_1 \in \Theta(\psi_1 - \psi_2)\}$, а количество дубликатов

$Occ(s', \psi') = Occ(s_1, \psi_1 - \psi_2)$, $s' \in \Theta(\psi')$ и $s' = s_1 \cup s_{R_2 \setminus R_1, NULL}$ (очевидно, что здесь и далее аналогичные представления строк единственны), а также

$$\left\langle \psi_2 - \psi_1, R_2 \right\rangle_{\varphi, R_2, R_1 \setminus R_2} \otimes \left\langle \{s_{R_1 \setminus R_2, NULL}^1\}, R_1 \setminus R_2 \right\rangle = \langle \psi'', R_1 \cup R_2 \rangle,$$

где основой мультимножества ψ'' является множество строк $\Theta(\psi'') = \{s_{R_1 \setminus R_2, NULL} \cup s_2 \mid s_2 \in \Theta(\psi_2 - \psi_1)\}$, а количество дубликатов

$Occ(s'', \psi'') = Occ(s_2, \psi_2 - \psi_1)$, $s'' \in \Theta(\psi'')$ и $s'' = s_{R_1 \setminus R_2, NULL} \cup s_2$.

Выше верхний индекс 1 в записи $\left\langle \{s_{R_2 \setminus R_1, NULL}^1\}, R_2 \setminus R_1 \right\rangle$ указывает на то, что строка $s_{R_2 \setminus R_1, NULL}$ входит в исходную таблицу только один раз, т.е. $\left\langle \{s_{R_2 \setminus R_1, NULL}^1\}, R_2 \setminus R_1 \right\rangle$ – константная таблица схемы $R_2 \setminus R_1$, первая компонента которой $\{1\}$ -мультимножество с основой $\{s_{R_2 \setminus R_1, NULL}\}$. Для таблицы $\left\langle \{s_{R_1 \setminus R_2, NULL}^1\}, R_1 \setminus R_2 \right\rangle$ полностью аналогично.

Под внешним левым соединением (Outer Left Join), индуцированным операцией φ , понимаем частичную бинарную операцию вида $\varphi_l : \Psi(R_1) \times \Psi(R_2) \rightarrow \Psi(R_1 \cup R_2)$, где $\text{dom } \varphi_l = \text{dom } \varphi$, $\varphi_l(\langle \psi_1, R_1 \rangle, \langle \psi_2, R_2 \rangle) = \varphi(\langle \psi_1, R_1 \rangle, \langle \psi_2, R_2 \rangle) \cup_{All}^{R_1 \cup R_2} \left\langle \psi_1 - \psi_2, R_1 \right\rangle_{\varphi, R_1, R_2 \setminus R_1} \otimes \left\langle \{s_{R_2 \setminus R_1, NULL}^1\}, R_2 \setminus R_1 \right\rangle$.

Под внешним правым соединением (Outer Right Join), индуцированным операцией φ , понимаем частичную бинарную операцию вида

$$\varphi_r : \Psi(R_1) \times \Psi(R_2) \xrightarrow{\sim} \Psi(R_1 \cup R_2), \quad \text{где } \text{dom } \varphi_r = \text{dom } \varphi, \quad \varphi_r(\langle \psi_1, R_1 \rangle, \langle \psi_2, R_2 \rangle) = \\ = \varphi(\langle \psi_1, R_1 \rangle, \langle \psi_2, R_2 \rangle) \cup_{All}^{R_1 \cup R_2} \left\langle \psi_2 - \psi_1, R_2 \right\rangle_{R_2, R_1 \setminus R_2} \otimes \left\langle \{s_{R_1 \setminus R_2, NULL}^1\}, R_1 \setminus R_2 \right\rangle.$$

Под внешним полным соединением (Outer Full Join), индуцированным операцией φ , понимаем частичную бинарную операцию вида $\varphi_f : \Psi(R_1) \times \Psi(R_2) \xrightarrow{\sim} \Psi(R_1 \cup R_2)$, где $\text{dom } \varphi_f = \text{dom } \varphi$,

$$\varphi_f(\langle \psi_1, R_1 \rangle, \langle \psi_2, R_2 \rangle) = \varphi(\langle \psi_1, R_1 \rangle, \langle \psi_2, R_2 \rangle) \cup_{All}^{R_1 \cup R_2} \left\langle \psi_1 - \psi_2, R_1 \right\rangle_{R_1, R_2 \setminus R_1} \otimes \\ \otimes_{R_1, R_2 \setminus R_1} \left\langle \{s_{R_2 \setminus R_1, NULL}^1\}, R_2 \setminus R_1 \right\rangle \cup_{All}^{R_1 \cup R_2} \left\langle \psi_2 - \psi_1, R_2 \right\rangle_{R_2, R_1 \setminus R_2} \otimes \\ \otimes_{R_2, R_1 \setminus R_2} \left\langle \{s_{R_1 \setminus R_2, NULL}^1\}, R_1 \setminus R_2 \right\rangle.$$

Под внешним соединением объединением (Outer Union Join), индуцированным операцией φ , понимаем частичную бинарную операцию вида $\varphi_U : \Psi(R_1) \times \Psi(R_2) \xrightarrow{\sim} \Psi(R_1 \cup R_2)$, где $\text{dom } \varphi_U = \text{dom } \varphi$, $\varphi_U(\langle \psi_1, R_1 \rangle, \langle \psi_2, R_2 \rangle) = \\ = \left\langle \psi_1 - \psi_2, R_1 \right\rangle_{R_1, R_2 \setminus R_1} \otimes \left\langle \{s_{R_2 \setminus R_1, NULL}^1\}, R_2 \setminus R_1 \right\rangle \cup_{All}^{R_1 \cup R_2} \left\langle \psi_2 - \psi_1, R_2 \right\rangle_{R_2, R_1 \setminus R_2} \otimes \\ \otimes_{R_2, R_1 \setminus R_2} \left\langle \{s_{R_1 \setminus R_2, NULL}^1\}, R_1 \setminus R_2 \right\rangle.$

6. Агрегатные операции

Широко используемыми (параметрическими) агрегатными операциями являются *Sum*, *Avg*, *Min*, *Max*, *Count*. Их аргументы – это конечные таблицы, а значения – одноатрибутные таблицы с одной строкой. Так, операция *Sum* рассчитывает сумму значений в соответствующем столбце заданной таблицы, при этом значения *NULL* игнорируются. Операция *Avg* определяет среднее арифметическое значений в соответствующем столбце заданной таблицы, при этом значения *NULL* игнорируются. Операции *Min* и *Max* находят наименьшее и наибольшее значения в соответствующем столбце заданной таблицы, при этом значения *NULL* также игнорируются. Операция *Count* определяет количество значений, отличных от *NULL*, в соответствующем столбце заданной таблицы. Операция *Count(*)* определяет количество строк в заданной таблице.

Пусть $\langle \psi, R \rangle \in \Psi(R)$, где ψ – конечное мультимножество и $A \in R$. Обозначим как α_A – мультимножество, которое содержит все элементы с учетом дубликатов столбца с атрибутом A таблицы $\langle \psi, R \rangle$. Тогда $\Theta(\alpha_A) = \{d \mid \exists s (s \in \Theta(\psi) \wedge \langle A, d \rangle \in s)\} = \{d \mid \{\langle A, d \rangle\} \in \Theta(\left(\pi_{\{A\}, R}(\langle \psi, R \rangle)\right)_1)\}$ – аналог

активного домена атрибута A относительно таблицы [5, 9]. Количество дубликатов элемента основы $d \in \Theta(\alpha_A)$ в мультимножестве α_A определяется как

$$\alpha_A(d) = \text{Occ}(\{\langle A, d \rangle\}, (\pi_{\{A\}, R}(\langle \psi, R \rangle))_1) = \sum_{\substack{s \in \Theta(\psi), \\ s(A)=d}} \text{Occ}(s, \psi). \quad \text{Пусть}$$

$2_m^{D'} = \{\alpha \mid \Theta(\alpha) \in 2^{D'}\}$ – семейство всех мультимножеств, основы которых являются конечными подмножествами множества D' ; здесь $D' \subseteq D$ – подмножество универсального домена.

Пусть Num – числовое подмножество универсального домена D , замкнутое относительно сложения. Множество Num расширим включением особого элемента $NULL$, но операцию сложения на случай, когда хотя бы один из аргументов является $NULL$, расширять не будем.

Зададим агрегатные операции Sum , Avg , Min , Max , $Count$. Общая схема: на конечном мультимножестве определяются функции суммирования, взятия наименьшего и наибольшего значений, определения среднего арифметического и количества элементов, а затем эти функции переносятся на таблицы. Заметим, что функции суммирования и нахождения среднего арифметического определены на конечном числовом мультимножестве.

Под операцией агрегирования $Sum_{A,R}$ по атрибуту A (конечных) таблиц схемы R , $A \in R$, понимаем унарную параметрическую операцию вида $Sum_{A,R} : \Psi(R) \rightarrow \Psi(\{A\})$, $Sum_{A,R}(\langle \psi, R \rangle) = \left\langle \left\{ \langle A, Sum(\alpha_A) \rangle \right\}^1, \{A\} \right\rangle$, где $\langle \psi, R \rangle \in \Psi(R)$, а Sum – функция, которая возвращает сумму значений столбца с атрибутом A таблицы $\langle \psi, R \rangle$ (эти значения могут повторяться), которые отличаются от значения $NULL$, кроме того предполагается, что этот столбец содержит только данные числового типа. Следовательно, $Sum : 2_m^{Num} \rightarrow Num$,

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

Как и раньше верхний индекс 1 в записи $\left\langle \left\{ \langle A, Sum(\alpha_A) \rangle \right\}^1, \{A\} \right\rangle$ указывает на то, что строка $\{\langle A, Sum(\alpha_A) \rangle\}$ входит в исходную таблицу только один раз, т.е. $\{\{\langle A, Sum(\alpha_A) \rangle\}^1\} - \{1\}$ -мультимножество.

Таким образом, имеем $Sum(\emptyset_m) = NULL$, $Sum(\{NULL^n\}) = NULL$, $Sum(\langle d_1^{n_1}, \dots, d_k^{n_k} \rangle) = \sum_{i=1}^k d_i n_i$, в предположении, что все элементы d_i отличны от элемента $NULL$.

Для случая пустой таблицы $\langle \psi_{\emptyset}, R \rangle$ операция агрегирования $Sum_{A,R}$ применяется так: $Sum_{A,R}(\langle \psi_{\emptyset}, R \rangle) = \left\langle \left\{ \left\langle A, NULL \right\rangle \right\}^{\uparrow}, \{A\} \right\rangle$, здесь $\psi_{\emptyset} = \emptyset_m$.

Пусть \leq – линейный порядок на универсальном домене \mathbf{D} . Под операцией агрегирования $Min_{A,R}$ по атрибуту A (конечных) таблиц схемы R , $A \in R$, понимаем унарную параметрическую операцию вида $Min_{A,R} : \Psi(R) \rightarrow \Psi(\{A\})$, $Min_{A,R}(\langle \psi, R \rangle) = \left\langle \left\{ \left\langle A, Min(\alpha_A) \right\rangle \right\}^{\uparrow}, \{A\} \right\rangle$, где $\langle \psi, R \rangle \in \Psi(R)$, а Min – функция, которая возвращает наименьшее значение среди значений столбца с атрибутом A таблицы $\langle \psi, R \rangle$, которые отличаются от значения $NULL$, т.е. $Min : 2_m^{\mathbf{D}} \rightarrow \mathbf{D}$,

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

Таким образом, имеем $Min(\emptyset_m) = NULL$, $Min(\{NULL^n\}) = NULL$, $Min(\{d_1^{n_1}, \dots, d_k^{n_k}\}) = \min\{d_1, \dots, d_k\}$, в предположении, что все элементы d_i , $i = \overline{1, k}$, отличны от элемента $NULL$.

Для случая пустой таблицы $\langle \psi_{\emptyset}, R \rangle$ операция агрегирования $Min_{A,R}$ применяется так: $Min_{A,R}(\langle \psi_{\emptyset}, R \rangle) = \left\langle \left\{ \left\langle A, NULL \right\rangle \right\}^{\uparrow}, \{A\} \right\rangle$.

Под операцией агрегирования $Max_{A,R}$ по атрибуту A (конечных) таблиц схемы R , $A \in R$, понимаем унарную параметрическую операцию вида $Max_{A,R} : \Psi(R) \rightarrow \Psi(\{A\})$, $Max_{A,R}(\langle \psi, R \rangle) = \left\langle \left\{ \left\langle A, Max(\alpha_A) \right\rangle \right\}^{\uparrow}, \{A\} \right\rangle$, где $\langle \psi, R \rangle \in \Psi(R)$, а Max – функция, которая возвращает наибольшее значение среди значений столбца с атрибутом A таблицы $\langle \psi, R \rangle$, которые отличаются от $NULL$, т.е. $Max : 2_m^{\mathbf{D}} \rightarrow \mathbf{D}$,

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

Таким образом, имеем $Max(\emptyset_m) = NULL$, $Max(\{NULL^n\}) = NULL$, $Max(\{d_1^{n_1}, \dots, d_k^{n_k}\}) = \max\{d_1, \dots, d_k\}$, в предположении, что все элементы d_i , $i = \overline{1, k}$, отличны от элемента $NULL$.

Для случая пустой таблицы $\langle \psi_{\emptyset}, R \rangle$ операция агрегирования $Max_{A,R}$ применяется так: $Max_{A,R}(\langle \psi_{\emptyset}, R \rangle) = \left\langle \left\{ \left\{ \langle A, NULL \rangle \right\} \right\}, \{A\} \right\rangle$.

Отметим, что функции Min и Max определяют наименьший или наибольший элементы основы мультимножества, которые отличаются от значения $NULL$, поэтому сравнимость особого элемента с остальными элементами универсального домена в данном случае несущественна. Это свойство важно при задании семантики фразы ORDER BY оператора запросов SELECT, отвечающей за упорядочение строк.

Под операцией агрегирования $Count_{A,R}$ по атрибуту A (конечных) таблиц схемы R , $A \in R$, понимаем унарную параметрическую операцию вида $Count_{A,R} : \Psi(R) \rightarrow \Psi(\{A\})$, $Count_{A,R}(\langle \psi, R \rangle) = \left\langle \left\{ \left\{ \langle A, Count(\alpha_A) \rangle \right\} \right\}, \{A\} \right\rangle$, где $\langle \psi, R \rangle \in \Psi(R)$, а $Count$ – функция, которая возвращает количество значений, которые отличаются от значения $NULL$, с учетом дубликатов в столбце с атрибутом A таблицы $\langle \psi, R \rangle$, т.е. $Count : 2_m^D \rightarrow N$, $Count(\alpha_A) = \sum_{d \in \Theta(\alpha_A) \setminus \{NULL\}}$

полагается по определению (и это естественно), что сумма пустого множества слагаемых равна нулю.

Таким образом, имеем $Count(\emptyset_m) = 0$, $Count(\{NULL^n\}) = 0$, $Count(\{d_1^{n_1}, \dots, d_k^{n_k}\}) = n_1 + \dots + n_k$, в предположении, что все элементы d_i , $i = \overline{1, k}$, отличны от элемента $NULL$.

Для случая пустой таблицы $\langle \psi_{\emptyset}, R \rangle$ операция агрегирования $Count_{A,R}$ применяется так: $Count_{A,R}(\langle \psi_{\emptyset}, R \rangle) = \left\langle \left\{ \left\{ \langle A, 0 \rangle \right\} \right\}, \{A\} \right\rangle$.

Допустим, что числовое подмножество Num универсального домена замкнуто относительно (частичной операции) деления $/: Num \times Num \xrightarrow{\sim} Num$. Доопределим операцию деления так, что когда первый аргумент равен $NULL$, то функция принимает значение $NULL$. Это связано с тем, что мы будем осуществлять суперпозиции и вместо первого аргумента подставлять значение функции Sum , а вместо второго – значение функции $Count$, учитывая, что функция $Count$ не может дать $NULL$.

Под операцией агрегирования $Avg_{A,R}$ по атрибуту A (конечных) таблиц схемы R , $A \in R$, понимаем унарную параметрическую операцию вида $Avg_{A,R} : \Psi(R) \rightarrow \Psi(\{A\})$, $Avg_{A,R}(\langle \psi, R \rangle) = \left\langle \left\{ \left\{ \langle A, Avg(\alpha_A) \rangle \right\} \right\}, \{A\} \right\rangle$, где $\langle \psi, R \rangle \in \Psi(R)$, а Avg – функция, которая возвращает среднее арифметическое значение элементов столбца с атрибутом A таблицы $\langle \psi, R \rangle$, которые

отличаются от значения $NULL$, с учетом дубликатов, т.е. $Avg : 2_m^{Num} \rightarrow Num$, $Avg(\alpha_A) = Sum(\alpha_A) / Count(\alpha_A)$.

Таким образом, из определений следуют равенства $Avg(\emptyset_m) = Sum(\emptyset_m) / Count(\emptyset_m) = NULL / 0 = NULL$,

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

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

предположении, что все элементы d_i , $i = \overline{1, k}$, отличны от элемента $NULL$.

Для случая пустой таблицы $\langle \psi_{\emptyset}, R \rangle$ операция агрегирования $Avg_{A,R}$ применяется так: $Avg_{A,R}(\langle \psi_{\emptyset}, R \rangle) = \langle \{ \langle A, NULL \rangle \}, \{A\} \rangle$.

Под операцией агрегирования $Count_{A,R}(\ast)$ (конечных) таблиц схемы R понимаем унарную параметрическую операцию вида

$Count_{A,R}(\ast) : \Psi(R) \rightarrow \Psi(\{A\})$, $Count_{A,R}(\ast)(\langle \psi, R \rangle) = \langle \{ \langle A, \|\psi\| \rangle \}, \{A\} \rangle$, где

$\langle \psi, R \rangle \in \Psi(R)$, а $\|\psi\|$ – это введенный ранее ранг мультимножества ψ .

Для случая пустой таблицы $\langle \psi_{\emptyset}, R \rangle$ операция агрегирования $Count_{A,R}(\ast)$ применяется так: $Count_{A,R}(\ast)(\langle \psi_{\emptyset}, R \rangle) = \langle \{ \langle A, \|\emptyset_m\| \rangle \}, \{A\} \rangle = \langle \{ \langle A, 0 \rangle \}, \{A\} \rangle$.

7. Выводы по результатам и направления дальнейших исследований

Существует ряд прикладных задач, особенностью которых является множественность и повторяемость данных. Например, социологические опросы различных групп населения, вычисления на ДНК, уточнения таблиц с дубликатами строк и другие. Математическим уточнением совокупностей с повторениями выступают мультимножества. Естественно возникает потребность в расширении возможностей реляционных баз данных за счет использования мультимножеств.

В данной работе рассматривается мультимножественная табличная алгебра, сигнатура которой пополнена новыми операциями: операциям внутренних и внешних соединений, операцией полусоединения, агрегатными операциями. Для определения внешних операций введен особый элемент универсального домена $NULL$.

Следует также отметить, что параметром агрегатной операции может выступать не только отдельный атрибут, но и некоторая функция над строкой. Соответствующее уточнение не трудно провести на основе представленных в работе построений.

В дальнейшем планируется исследовать свойства операций мультимножественной табличной алгебры.

ЛИТЕРАТУРА

1. Codd E.F. Relational Completeness of Data Base Sublanguages / Codd E.F. // Data Base Systems. – New York: Prentice-Hall. – 1972. – P. 65-93.
2. Lamperti G. On Multisets in Database Systems / G. Lamperti, M. Melchiori, M. Zanella // 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.
3. Гарсиа-Молина Г. Системы баз данных: [полный курс: пер. с англ.] / Г. Гарсиа-Молина, Дж. Ульман, Дж. Уидом. – Москва: "Вильямс", 2004. – 1088 с.
4. Silbeschatz A., Korth H., Sudarshan S. Database System Concepts. – McGraw-Hill, 2011. – 1376 p.
5. Реляційні бази даних: табличні алгебри та SQL-подібні мови / В.Н. Редько, Ю.Й. Брона, Д.Б. Буй, С.А. Поляков. – Київ: Видавничий дім "Академперіодика", 2001. – 198 с.
6. Буй Д.Б. Сучасний стан теорії мультимножин / Д.Б. Буй, Ю.О. Богатирьова // Вісник Київського університету. Сер.: фіз.-мат. науки. – 2010. – Вип. 1. – С. 51-58.
7. Петровский А.Б. Пространства множеств и мультимножеств / А.Б. Петровский. – Москва: "Едиториал УРСС", 2003. – 248 с.
8. Глушко І.М. Мультимножинна таблична алгебра / І.М. Глушко // Proceedings of the International Scientific Conference of Student and Young Scientists "Theoretical and Applied Aspects of Cybernetics" (ТААС'2011, Kyiv, February 21–25, 2011). – P. 77-79.
9. Мейер Д. Теория реляционных баз данных: [пер. с англ.] / Д. Мейер. – Москва: Мир, 1987. – 608 с.