

## ИНФОРМАТИКА И ПРОГРАММИРОВАНИЕ

УДК 004.652

О.А. Бистерфельд

### ОЦЕНКА РЕСУРСОВ ПАМЯТИ, НЕОБХОДИМЫХ ДЛЯ РЕАЛИЗАЦИИ КАТЕГОРИРОВАННЫХ ОТНОШЕНИЙ В РЕЛЯЦИОННЫХ БАЗАХ ДАННЫХ

Предложена модель базы категорированных данных, учитывающая объем и особенности распределения данных по категориям. Выведены формулы зависимости требуемых ресурсов памяти от параметров структуры категорированных данных для наиболее распространенных вариантов представления категорированных отношений. Аналитические выражения метода уточнены с помощью имитационной программы. Получаемые оценки используются при решении задачи выбора вариантов представления категорированных данных.

**Ключевые слова:** *базы данных, отношения категоризации, объем памяти.*

#### 1. Способы представления категорированных отношений

В базах данных (БД) информационных систем значительная часть данных связана отношениями категоризации.

Наиболее развиты способы реализации категорированных отношений в СУБД Oracle. Варианты СУБД Oracle охватывают практически все возможные варианты других СУБД, поэтому они приняты за основу при разработке аналитического метода оценки объемов памяти, необходимых для реализации категорированных отношений в реляционных БД. В многотомном описании методологии проектирования Oracle [1, с. 71–76.] приведены только словесные описания вариантов реализации с кратким перечнем достоинств и недостатков каждого.

На рис. 1 показана модель отношения категоризации по нотации Баркера («супертип» – «тип»).

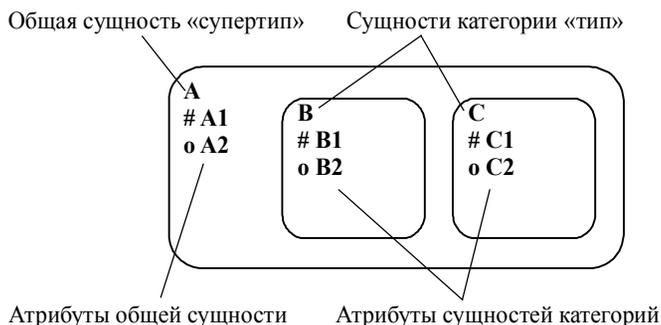


Рис. 1. Базовая ER-диаграмма

Категорированные данные могут быть представлены (рис. 2):

- в одной таблице (вариант *a*);
- в нескольких (по количеству категорий) таблицах (вариант *б*);
- в нескольких (по количеству категорий + одна) таблицах (варианты *в*, *г*).

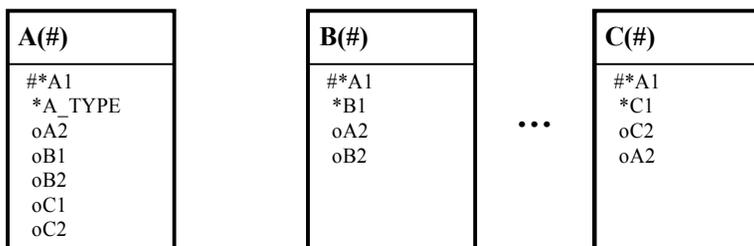
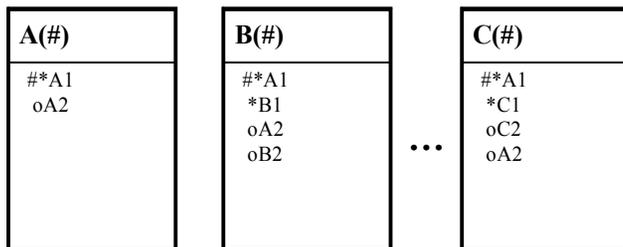
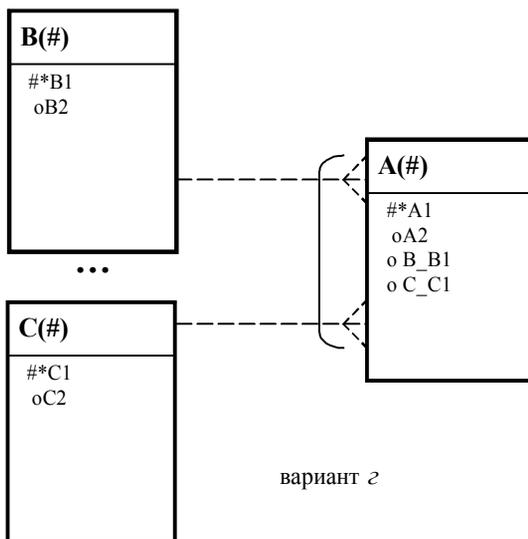
вариант *a*вариант *б*вариант *в*вариант *г*

Рис. 2. Варианты реализации категорированных отношений

## 2. Аналитический метод оценки различий в требуемых ресурсах памяти БД вариантов представления категорированных отношений

Введем обозначения:  $N$  – число записей в БД;  $N_{\max}$  – максимально возможное число записей в БД;  $K_{y,и}$  – размер уникального идентификатора;  $A_{o,a}$  – размер общих атрибутов;  $K_k$  – число категорий;  $A_i$  – размер частных атрибутов (для категории  $i$ );  $N_i$  – число записей категории  $i$ ;  $N_{i\max}$  – максимально возможное число записей категории  $i$ ;  $K_{y,иi}$  – размер уникального идентификатора для категории  $i$ . Должны выполняться следующие неравенства:

$$K_{y,и} \geq \log_2 N_{\max}; K_{y,иi} \geq \log_2 N_{i\max}, \quad (1)$$

однако на практике для уникальных идентификаторов записей используют десятичные числа, поэтому имеем

$$K_{y,и} \geq \lg N_{\max}; K_{y,иi} \geq \lg N_{i\max}. \quad (2)$$

Кроме того,

$$\sum_{i=1}^{k_k} N_i = N. \quad (3)$$

Пределы суммирования во всех формулах данной статьи одинаковы, поэтому ниже они не указываются (вместо  $\sum_{i=1}^{k_k}$  используется  $\Sigma$ ).

В некоторых упрощениях (при несущественном влиянии на результат) принимается также, что  $K_{y,иi} = K_{y,и}$ .

Объем памяти по варианту  $a$  –  $O_a$  (без учета несущественного по объему атрибута TYPE):

$$O_a = N(K_{y,и} + A_{o,a} + \Sigma A_i). \quad (4)$$

Объем по варианту  $b$  – ( $O_b$ ):

$$O_b = \Sigma N_i(K_{y,и} + A_{o,a} + A_i). \quad (5)$$

Объем по варианту  $в$  – ( $O_в$ ):

$$O_в = N(K_{y,и} + A_{o,a}) + \Sigma N_i(K_{y,и} + A_{o,a} + A_i). \quad (6)$$

Объем по варианту  $г$  – ( $O_г$ ):

$$O_г = N(K_{y,и} + A_{o,a} + \Sigma K_{y,иi}) + \Sigma N_i(K_{y,иi} + A_i). \quad (7)$$

Для анализа эффективности различных вариантов представления категорированных отношений в реляционных БД представляют интерес:

- относительные оценки:  $(O_a - O_b)/O_a$ ;  $(O_a - O_в)/O_a$ ;  $(O_a - O_г)/O_a$ ;

- абсолютные оценки:  $O_a$ ;  $O_a - O_b$ ;  $O_a - O_в$ ;  $O_a - O_г$ .

С учетом формул (1) – (7):

$$(O_a - O_b)/O_a = \{N(K_{y,и} + A_{o,a} + \Sigma A_i) - \Sigma N_i(K_{y,и} + A_{o,a} + A_i)\} / N(K_{y,и} + A_{o,a} + \Sigma A_i); \quad (8)$$

$$(O_a - O_в)/O_a = \{N(K_{y,и} + A_{o,a} + \Sigma A_i) - [N(K_{y,и} + A_{o,a}) + \Sigma N_i(K_{y,и} + A_{o,a} + A_i)]\} / N(K_{y,и} + A_{o,a} + \Sigma A_i); \quad (9)$$

$$(O_a - O_г)/O_a = \{N(K_{y,и} + A_{o,a} + \Sigma A_i) - [N(K_{y,и} + A_{o,a} + \Sigma K_{y,иi}) + \Sigma N_i(K_{y,иi} + A_i)]\} / N(K_{y,и} + A_{o,a} + \Sigma A_i). \quad (10)$$

После преобразований, с учетом формул (2), (3) и некоторых приближений:

$$(O_a - O_b)/O_a = A_{o,a} [\Sigma A_i / A_{o,a} - \Sigma (N_i / N) (A_i / A_{o,a})] / [A_{o,a} (1 + \Sigma A_i / A_{o,a}) + \lg N_{\max}]; \quad (11)$$

$$O_a - O_в = N A_{o,a} [\Sigma A_i / A_{o,a} - \Sigma (N_i / N) (A_i / A_{o,a})]; \quad (12)$$

$$O_a - O_{\bar{e}}/O_a = \{A_{o,a}[\Sigma A_i/A_{o,a} - 1 - \Sigma(N_i/N)(A_i/A_{o,a})] - \lg N_{\max}\} / [A_{o,a}(1 + \Sigma A_i/A_{o,a}) + \lg N_{\max}]; \quad (13)$$

$$O_a - O_{\bar{e}} = N\{A_{o,a}[\Sigma A_i/A_{o,a} - 1 - \Sigma(N_i/N)(A_i/A_{o,a})] - \lg N_{\max}\}; \quad (14)$$

$$(O_a - O_z)/O_a = \{A_{o,a}[\Sigma A_i/A_{o,a} - \Sigma(N_i/N)(A_i/A_{o,a})] - \Sigma \lg N_{i\max} - \lg N_{\max}\} / [A_{o,a}(1 + \Sigma A_i/A_{o,a}) + \lg N_{\max}]; \quad (15)$$

$$(O_a - O_z) = N\{A_{o,a}[\Sigma A_i/A_{o,a} - \Sigma(N_i/N)(A_i/A_{o,a})] - \Sigma \lg N_{i\max} - \lg N_{\max}\}; \quad (16)$$

$$O_a = N[A_{o,a}(1 + \Sigma A_i/A_{o,a}) + \lg N_{\max}]; \quad (17)$$

$$O_{\bar{e}} = N[A_{o,a}\{1 + \Sigma(N_i/N)(A_i/A_{o,a})\} + \lg N_{\max}]; \quad (18)$$

$$O_{\bar{e}} = N[A_{o,a}\{2 + \Sigma(N_i/N)(A_i/A_{o,a})\} + 2\lg N_{\max}]; \quad (19)$$

$$O_z = N[A_{o,a}\{1 + \Sigma(N_i/N)(A_i/A_{o,a})\} + \Sigma \lg N_{i\max} + 2\lg N_{\max}]. \quad (20)$$

Выражения (11) – (16) могут быть использованы для проведения исследований и анализа эффективности различных вариантов представления категорированных отношений в реляционных БД [2, с. 86–87].

### 3. Модель базы категорированных данных

Следует обратить внимание на набор компонентов выражений (11) – (20):

$$\{N; N_{\max}; N_{i\max}; A_{o,a}; \Sigma A_i/A_{o,a}; \Sigma(N_i/N)(A_i/A_{o,a})\}.$$

Эти компоненты достаточно адекватно представляют реальную БД и особенности структуры данных.

Введем обозначения:  $a_i = A_i/A_{o,a}$ ;  $n_i = N_i/N$  и преобразуем формулы (11) – (20):

$$(O_a - O_{\bar{e}})/O_a = A_{o,a}(\Sigma a_i - \Sigma n_i a_i) / [A_{o,a}(1 + \Sigma a_i) + \lg N_{\max}]; \quad (21)$$

$$O_a - O_{\bar{e}} = NA_{o,a}(\Sigma a_i - \Sigma n_i a_i); \quad (22)$$

$$(O_a - O_{\bar{e}})/O_a = [A_{o,a}(\Sigma a_i - 1 - \Sigma n_i a_i) - \lg N_{\max}] / [A_{o,a}(1 + \Sigma a_i) + \lg N_{\max}]; \quad (23)$$

$$O_a - O_{\bar{e}} = N[A_{o,a}(\Sigma a_i - 1 - \Sigma n_i a_i) - \lg N_{\max}]; \quad (24)$$

$$(O_a - O_z)/O_a = [A_{o,a}(\Sigma a_i - \Sigma n_i a_i) - \Sigma \lg N_{i\max} - \lg N_{\max}] / [A_{o,a}(1 + \Sigma a_i) + \lg N_{\max}]; \quad (25)$$

$$(O_a - O_z) = N[A_{o,a}(\Sigma a_i - \Sigma n_i a_i) - \Sigma \lg N_{i\max} - \lg N_{\max}]; \quad (26)$$

$$O_a = N[A_{o,a}(1 + \Sigma a_i) + \lg N_{\max}]; \quad (27)$$

$$O_{\bar{e}} = N[A_{o,a}(1 + \Sigma n_i a_i) + \lg N_{\max}]; \quad (28)$$

$$O_{\bar{e}} = N[A_{o,a}(2 + \Sigma n_i a_i) + 2\lg N_{\max}]; \quad (29)$$

$$O_z = N[A_{o,a}(1 + \Sigma n_i a_i) + \Sigma \lg N_{i\max} + 2\lg N_{\max}]. \quad (30)$$

Формулы (21) – (30) положены в основу аналитического метода оценки требуемых объемов памяти по вариантам представления категорированных отношений.

Целесообразно использовать набор  $\{N; N_{\max}; |N_{i\max}|; A_{o,a}; |a_i|; |n_i|\}$  как модель базы категорированных данных для исследований зависимостей объемов памяти, необходимых для реализации категорированных отношений.

Замечание. Если  $\Sigma a_i$  имеет простой «физический» смысл (эта сумма показывает соотношение между суммарными размерами общих и частных атрибутов), то с компонентом  $\Sigma n_i a_i$  несколько сложнее. Диапазон значений  $\Sigma n_i a_i$ :  $-0 < \Sigma n_i a_i < \Sigma a_i$ . Если значения  $\Sigma n_i a_i$  в левой части диапазона (от  $0,5 \Sigma a_i$  к 0), то в БД незначительна доля записей, относящихся к категориям, имеющим существенные размеры част-

ных атрибутов, и, наоборот, существенна доля записей, относящихся к категориям, имеющим незначительные размеры частных атрибутов. Если значения  $\sum n_i a_i$  в правой части диапазона (от  $0,5 \Sigma a_i$  к  $\Sigma a_i$ ), то в БД значительна доля записей, относящихся к категориям, имеющим существенные размеры частных атрибутов, и, наоборот, незначительна доля записей, относящихся к категориям, имеющим незначительные размеры частных атрибутов.

#### 4. Зависимости затрат памяти от параметров модели базы категорированных данных

Некоторые зависимости затрат памяти, полученные по выражениям (21)-(30), приведены на диаграммах (рис. 3, 4).

Зависимость затрат памяти от структуры категорированных записей ( $\Sigma a_i n_i$ )  
( $N = 500\ 000$ ;  $\Sigma a_i = 10$ ;  $\lg N_{\max} = 6$ ;  $A_{o.a} = 512$  симв.)

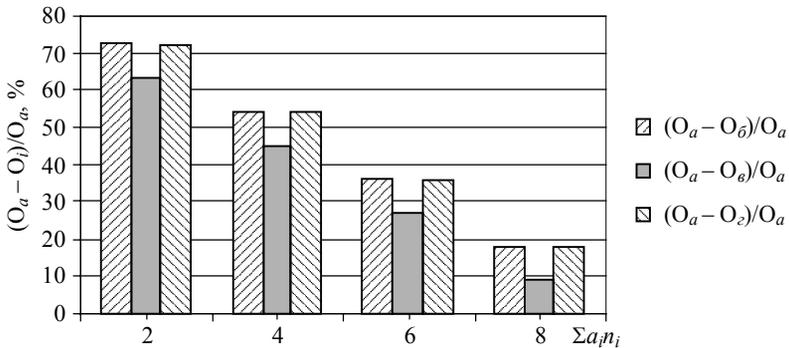


Рис. 3. Уменьшение (экономия) ресурсов памяти при переходе от варианта представления категорированных данных  $a$  к вариантам  $b, в, з$  в зависимости от структуры категорированных данных ( $\Sigma a_i n_i$ ) в БД

Зависимость затрат памяти от соотношения общих и частных атрибутов ( $\Sigma a_i$ )  
( $\Sigma n_i a_i = 0,5 \Sigma a_i$ ;  $N = 500\ 000$ ;  $\lg N_{\max} = 6$ ;  $A_{o.a} = 512$  симв.)

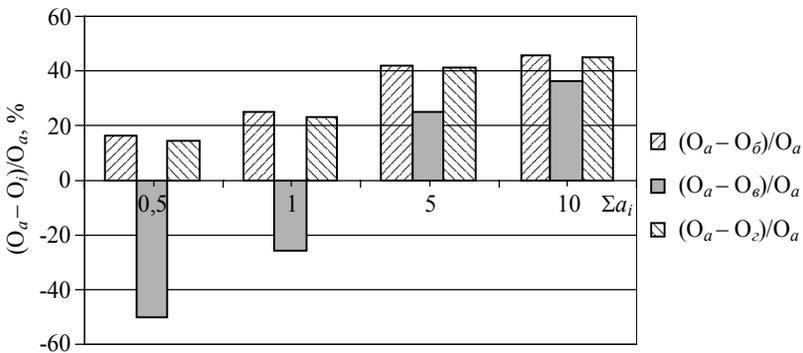


Рис. 4. Уменьшение (экономия) ресурсов памяти при переходе от варианта представления категорированных данных  $a$  к вариантам  $b, в и з$  в зависимости от соотношения общих и частных атрибутов ( $\Sigma a_i$ )

Полученные с помощью аналитических выражений данные показывают потенциальные возможности минимизации затрат ресурсов памяти с помощью выбора вариантов представления категорированных данных. В предельных случаях ресурсы памяти, необходимые для реализации категорированных данных, могут быть сокращены на 60 – 70 %.

### 5. Имитационная программа для уточнения аналитических выражений метода оценки требуемых ресурсов памяти

В аналитических выражениях метода оценки объемов памяти (разделы 2, 3) единицей измерения объема памяти является символ данных, заносимых в базу данных –  $S_d$ . Такой выбор удобен для определения параметров модели категорированных данных по анализу информационных моделей. В то же время при оценке затрат ресурсов памяти общепринятыми являются единицы измерения: бит, кбит, Мбит и Гбит.

Помимо этого, в ЭВМ и в СУБД используются методы сжатия данных. При записи таких типов данных, как числовые, несмотря на то, что отдельным атрибутам (в СУБД) резервируются определенные максимально возможные значения, нулевые значения слева от значащих символов числа не записываются и не сохраняются в памяти ЭВМ. Как правило, не заносятся (и не требуют ресурсов памяти) пустые значения атрибутов.

Современные СУБД обладают развитыми средствами повышения производительности баз данных. Ряд методов повышения производительности достигают за счет избыточного хранения данных (например, методы индексации записей в таблицах БД).

Для практического применения аналитических выражений (21) – (30) необходимо определить коэффициент перевода единиц измерения, учета сжатия и избыточности данных –  $k_{ес}$  размерностью бит/символ данных. При этом оценки (в единицах измерения – бит) необходимых затрат ресурсов памяти для реализации различных вариантов представления категорированных отношений:

$$O_a^{\text{бит}} = k_{ес}^a O_a; O_b^{\text{бит}} = k_{ес}^b O_b; O_e^{\text{бит}} = k_{ес}^e O_e; O_z^{\text{бит}} = k_{ес}^z O_z.$$

Для перевода единиц измерения необходимо учитывать особенности кодирования символов данных —  $S_d$ . Каждый символ данных представляется обычно десятичным или шестнадцатеричным кодом, содержащим несколько символов (символов кода –  $S_k$ ). Наиболее употребительным является стандарт ASCII (American Standard Code for Information Interchange, стандарт ANSI – американского национального института стандартов). Часто используемая для передачи данных американская версия семибитовой кодировки символов кода ( $S_k$ ) утверждена ISO. Восьмой бит символа кода ASCII обычно является битом контроля четности (паритета). Таким образом, одним из компонентов  $k_{ес}$  является коэффициент ( $n_{S_k}$ ) двоичного представления символа кода  $S_k$ . Для ASCII  $n_{S_k} = 8$  (бит/символ);  $n_{S_k} = 1$  (байт/символ);  $n_{S_k} = 1/1024$  (кбайт/символ) и т.д.

На ПК используется так называемый расширенный код ASCII, в котором первые 128 комбинаций совпадают со стандартным, а остальные используются для представления национальных алфавитов, псевдографики и специальных знаков. Кодировка символов данных отличается количеством знаков кода ( $N_{S_k}$ ). Например, цифры и большая часть латинских символов кодируются двумя десятичными знаками ( $N_{S_k} = 2$ ), а символы кириллицы — тремя ( $N_{S_k} = 3$ ). Необходимый для

представления данных объем памяти зависит от набора символов, необходимого для определенного типа данных, заносимых в базу данных. В первом приближении может быть принято, что набор символов данных носит случайный характер, и этот набор может быть представлен математическим ожиданием значения числа символов кода –  $\overline{N_{S_k}}$ . По стандарту ASCII для большинства наборов символов значение  $\overline{N_{S_k}}$  будет находиться в диапазоне от 2 до 3.

Таким же образом может быть признана случайной величина коэффициента учета сжатия ( $k_{сж}$ ) данных при записи в ЭВМ (при записи данных сжатие проводится только по отдельным атрибутам, и степень сжатия различна). Фактор сжатия данных должен приводить к уменьшению  $k_{сж}$ . При сжатии данных, например в 1,5 – 2 раза, общая оценка  $k_{сж}$  должна находиться в диапазоне от 1 до 2 (при измерении в единицах байт/символ).

Избыточность данных, прежде всего при использовании индексирования записей, также может быть учтена коэффициентом ( $k_{и.д}$ ). При индексировании записей СУБД может создавать дополнительные колонки, с помощью которых упорядочиваются записи. В связи с этим введение  $k_{и.д}$  адекватно отражает изменения в требуемых ресурсах памяти.

С учетом представления коэффициента сжатия средним значением общий вид  $k_{сж}$ :

$$k_{сж} = k_{сж} k_{и.д} n_{S_k} N_{S_k}.$$

Значения  $n_{S_k}$  и  $k_{и.д}$  могут быть определены при достаточно тщательном анализе конкретных вариантов реализаций фрагментов БД с категорированными данными. Аналитически определить  $k_{сж}$  и  $N_{S_k}$  гораздо сложнее. Предлагается использовать для оценки  $k_{сж}$  имитационную программу без детального анализа составляющих  $k_{сж}$ . Программа представляет собой тестовую базу данных со структурой, соответствующей исследуемой структуре категорированных данных (соотношение общих и частных атрибутов, в соответствии с моделью категорированных данных, раздел 3) и вариантам представления категорированных отношений.

Шаги процедуры определения значения  $k_{сж}$ :

- внесение набора числа тестовых записей (ряд значений  $N$ ) и фиксация затрат (в битах) ресурсов памяти –  $O_a^{имит}(N)$ ;  $O_b^{имит}(N)$ ;  $O_c^{имит}(N)$ ;  $O_z^{имит}(N)$ ;
- расчет затрат ресурсов по выражениям (27) – (30) в символах ( $S_n$ ) для каждого значения из ряда  $N$  –  $O_a(N)$ ;  $O_b(N)$ ;  $O_c(N)$ ;  $O_z(N)$ ;
- расчет  $k_{сж}$  для каждого  $N_i$ :

$$k_{сж}^a(N) = O_a^{имит}(N) / O_a(N); k_{сж}^b(N) = O_b^{имит}(N) / O_b(N);$$

$$k_{сж}^c(N) = O_c^{имит}(N) / O_c(N); k_{сж}^z(N) = O_z^{имит}(N) / O_z(N);$$

- определение среднего значения  $k_{сж}^a$ ;  $k_{сж}^b$ ;  $k_{сж}^c$ ;  $k_{сж}^z$ ;

- построение зависимостей  $O_a^{бит}$   $O_b^{бит}$   $O_c^{бит}$   $O_z^{бит}$ .

Использование оценки на основе среднего значения ( $k_{сж}$ ), получаемого по данным имитационной программы, позволяет получить аналитические выражения, отличающиеся от данных имитационной программы на 3 – 5 %, что вполне допустимо для использования их для оценок затрат ресурсов при проектировании реальных баз данных.

Величина  $k_{сж}(N)$  при различных значениях  $N$  и вариантах структуры категорированных данных находилась в пределах от 0,8 до 3,2 (оценка диапазона при анализе факторов, влияющих на  $k_{сж}(N)$ , – от 1 до 2). Расширение диапазона влево

можно объяснить тем, что при малом числе записей более существенно влияет фактор сжатия данных при записи в БД (например, для значений уникальных атрибутов достаточно одного-двух символов вместо  $\lg N_{\max}$ ). Расширение диапазона вправо до значения  $k_{\text{сж}}(N) = 3,2$  также объяснимо. Часть тестовых данных использовалась на текстовых фрагментах, содержащих только символы кириллицы, которые представляются тремя символами кода по ASCII. Кроме того, дополнительные затраты ресурсов памяти связаны с индексацией первичных ключей.

### Заключение

Оценить различия в объеме памяти БД вариантов представления категорированных отношений возможно аналитическими расчетами. Аналитические выражения определяют зависимости требуемых объемов памяти от параметров структуры категорированных данных для наиболее распространенных на практике вариантов представления категорированных отношений.

В аналитической модели должны быть учтены особенности структуры категорированных данных. Модель особенностей категорированных данных используется в аналитических выражениях метода оценки различий в объеме памяти БД вариантов представления категорированных отношений. Представление моделью реальных баз категорированных данных распространяет конкретные аналитические выражения метода на группу баз категорированных данных и сокращает затраты на моделирование.

Уточнить аналитические выражения метода оценки требуемых ресурсов памяти можно с помощью имитационной программы. Программа обеспечивает проверку достоверности аналитических выражений метода. Данные, получаемые имитационной программой, используются для определения коэффициента учета сжатия данных в СУБД, особенностей кодирования символов и дополнительных затрат ресурсов памяти при использовании в СУБД, например избыточных методов повышения производительности.

### ЛИТЕРАТУРА

1. *CDM* – метод разработки информационных систем фирмы Oracle // Oracle Magazine: russian edition. 1997. № 2.
2. Бистерфельд О.А., Сидоров М.В., Таганов Р.А. Исследование зависимости затрат памяти на представление категорированных отношений в реляционных базах данных // Новые информационные технологии в научных исследованиях и в образовании: тез. докл. 4-й Всероссийской научно-технической конференции. Рязань, 1999.

*Бистерфельд Ольга Александровна*

Рязанский государственный университет имени С.А. Есенина

E-mail: o.bisterfeld@rsu.edu.ru

Поступила в редакцию 18 июля 2011 г.