

МИНИСТЕРСТВО ОБРАЗОВАНИЯ *УКРАИНЫ
КИЕВСКИЙ ГОСУДАРСТВЕННЫЙ ЭКОНОМИЧЕСКИЙ УНИВЕРСИТЕТ

на правах рукописи

АМАНА НАХУДА ИБРАХИМ

ОПТИМИЗАЦИЯ СТРУКТУР БАЗ ДАННЫХ

Специальность 08.00.13. – Экономико-математические
методы

А В Т О Р Е Ф Е Р А Т
диссертация на соискание ученой степени
кандидата экономических наук

Киев 1993

0.4



00343926 (R)

Работа выполнена на ка
ковском инженерно-экономичес

Научный руководитель -

доктор технических наук, профессор
Волколупова Р.Т.

Официальные оппоненты:

Доктор экономических наук, профессор
Суслов О.П.

Кандидат экономических наук, доцент
Витлинский В.В.

Ведущая организация:

Главный научно-исследовательский ин-
ститут по проблемам информатики
Министерства экономики Украины

Защита диссертации состоится "24" июня 1993 года
в 14 часов на заседании Специализированного совета
К.068.28.05. в Киевском государственном экономическом университете
/252057, г.Киев-57, Проспект Победы, 54/1, ауд. 214 /.

С диссертацией можно ознакомиться в библиотеке университета.

Автореферат разослан "24" мая 1993 года.

Ученый секретарь
Специализированного совета

В.П.Кулагина

ЛНБ ім. В. Стефаніка
АН України

I. ОБЩАЯ ХАРАКТЕРИСТИКА

I.1. Актуальность темы. Одной из важнейших задач любой системы управления является переработка информации. Важную роль в переработке информации призвана сыграть технология создания баз данных. Важнейшим элементом этой технологии является процесс выбора структуры базы данных. Выбор надежного представления такой структуры по тому или иному критерию составляет задачу логического проектирования, решение которой позволит оптимизировать базу данных, что обеспечит эффективность функционирования информационной системы по времени и стоимости.

Постоянный поиск методов, позволяющих снизить затраты на стоимость хранения данных и их переработку, послужил объективной необходимостью углубления исследований по данным вопросам и определил актуальность темы диссертационной работы.

I.2. Цель и задачи исследования. Целью диссертационной работы является исследование, разработка методов, алгоритмов и программ решения задач оптимизации структур баз данных и построение оптимальной логической структуры.

Достижение поставленной цели обеспечивается решением следующих задач:

1. Создание модифицированного метода динамического сгущения.
2. Разработка метода и алгоритма исключения нулевых элементов.
3. Построение альтернативного множества.
4. Нахождение локального оптимума.
5. Создание модифицированного метода выбора системы отношений между элементами информации.

6. Выделение классов однотипных элементов.

7. Создание метода и алгоритма исключения бесперспективных вариантов.

8. Разработка метода и алгоритма синтеза оптимальной структуры базы данных.

1.3. Методы исследования. При выполнении данной работы применялись методы исследования операций, векторной оптимизации, кластер-анализа, теории множеств, теории матриц, теории графов.

1.4. Научная новизна результатов исследований:

- разработаны методы: модификация динамического сгущения, исключения нулевых элементов, модификация метода выбора системы отношений между элементами информации, исключения бесперспективных вариантов, векторной оптимизации;

- созданы алгоритмы и программы вышеперечисленных методов.

1.5. Практическая значимость. Разработанные методы и алгоритмы могут быть использованы для данных, позволяя оптимизировать структуры всех типов данных: иерархических, сетевых, реляционных и структуру представленную в виде гиперграфов.

1.6. Реализация результатов работы. Основные результаты работы использованы текстильным заводом COTEMA на Мадагаскаре для перспективного развития информационных систем (свидетельство посольства Республики Мадагаскар №318 AMBAMOSC/AFCULT).

1.7. Апробация работы. Достоверность полученных результатов и выводов подтверждена расчетом специальных математических критериев, апробацией методов оптимизации логической структуры баз данных на тестовых примерах, докладами на учредительной конференции международной ассоциации по нетрадиционным методам оптимизации (г. Дивногорск, 16 марта 1992 г.), в Междуна-

родной школе "Проектирование автоматизированных систем контроля и управления сложным объектом" (г.Туапсе 17 сентября 1992 г.), а также регулярно в течение всего периода выполнения исследований на заседаниях семинара "Проблемы экономической кибернетики" Академии наук Украины (г.Харьков, ХИЭИ, кафедре информационных систем), кроме того доклад по теме диссертации направлялся во Францию на конференцию: *10^{ème} conférence internationale sur l'analyse et l'optimisation des Systèmes approchés Fréquentielles et temporelles des Systèmes de Dimension infinie, 9-12 Juin 1992, SOPHIA - ANTIPOLIS (FRANCE)* и был принят экспертной комиссией.

1.8. Публикации. По теме диссертации опубликовано 3 печатных работы.

1.9. Структура и объем диссертации. Диссертационная работа состоит из введения, трех глав, выводов по главам, заключения, списка использованной литературы и трех приложений, 7 рисунков, 8 таблиц. Список литературы включает 124 наименования, из них 9 на иностранных языках.

2. СОДЕРЖАНИЕ РАБОТЫ

Во введении обосновываются актуальность темы, выбор предмета исследований, научная новизна и практическая значимость полученных результатов, формулируются основные направления исследований.

Первая глава посвящена критическому анализу некоторых существующих методов оптимизации логической структуры баз данных, а также сформулированы задачи исследования.

В первых шагах проектирования структур баз данных требуется определение состава информационных массивов. Традицион-

ный метод сортировки массивов требует большого перебора и более того, не устраняет дублирования информации, что приводит к неоптимальной логической структуре базы данных, а значит к ухудшению качества управления.

Информация, подлежащая обработке при решении задач, представляет собой информационный массив о некотором множестве элементов исследуемого объекта. Данные о каждом отдельном элементе такого множества составляют определенную совокупность, которую будем называть информационным модулем (ИМ). Информационный модуль расчленяется на отдельные части, называемые элементарными информационными модулями (ЭИМ) и представляющие собой последовательности символов определенного алфавита. Отдельные ЭИМ несут законченную смысловую нагрузку, соответствующую качественным и количественным характеристикам элементов объектов, описываемым информационным модулем. Тогда для устранения ранее перечисленных недостатков требуется образовать альтернативные модули наименьшей мощности из всех возможных скомпонованных по исходным данным ИМ и ввести оптимизацию этого альтернативного множества согласно соответствующим критериям.

В целях получения набора альтернативных информационных модулей предлагается разработать модифицированный метод динамического сгущения и метод исключения нулевых элементов, целью которых является устранение дублирования информации в покрытии и снижение затрат на вычисление при работе с разреженной матрицей. Выбор рациональных по составу и количеству альтернативных ИМ осуществляется с помощью метода нахождения локального оптимума.

Все эти предлагаемые разработанные методы направлены на повышение качества управления.

Рассмотренный метод приведения отношений к третьей норма-

льной форме затрудняет проектирование подходов базы данных, поскольку этот метод не учитывает наличие всех решаемых задач. В этой связи, для устранения этого недостатка необходимо применить модифицированный метод выбора системы отношений.

Оценка логической структуры базы данных осуществляется методом многокритериальной оптимизации, компромиссным путем, с разработанным методом исключения бесперспективных вариантов.

Во второй главе разработаны модификация метода динамического сгущения, метод исключения нулевых элементов, применен метод локального оптимума, разработан модифицированный метод выбора системы отношений, сформулирована математическая модель названных методов, применен подход к построению альтернативного множества и приведены результаты действия алгоритмов этих методов.

Введем некоторые понятия и обозначения.

Множества элементов, определяющее ИМ, необходимое для решения задач Z_j ($j = \overline{1, n}$) обозначим как $\Omega = \{P_1, P_2, \dots, P_l, \dots, P_k\}$.

Множество непустых подмножеств $P_l \subset \Omega$, ($l = 1, 2, \dots, k$), объединения которых заполняют Ω , назовем его покрытием Ω' со свойствами:

$$1. \forall l \in \{1, 2, \dots, k\}, l \leq k, P_l \neq \emptyset,$$

$$2. \bigcup_{l=1}^k P_l = \Omega'.$$

Тогда совокупность $P_l \subset \Omega'$ ($i \leq l$), попарное пересечение которых со всеми совокупностями P_j ($j \leq l, i \neq j$) образует пустые множества, определяет разбиение Ω' на классы R_i, P_j , предельное количество которых равно n , вычисление которого происходит с помощью соответствующего алгоритма.

Свойства классов:

$$1. P_i \cap P_j = \emptyset, \text{ где } P_{ij} = (P_3 \cup \dots \cup P_m) \cap (P_k \cup \dots \cup P_r),$$

2. $\forall s, m, k, \tau \in \{1, 2, \dots, l\}$.

Проектирование информационных модулей можно сформулировать как задачу:

- определения рациональных по составу модулей в том смысле, что каждый информационный модуль содержит минимальные избыточные и максимальные необходимые элементарные модули для решения задач;

- определения минимального количества набора информационных модулей.

Данная задача разделена на две подзадачи.

Первая подзадача: образование альтернативных ИМ наименьшей мощности, а вторая - оптимизация этого альтернативного множества.

Для решения первой задачи предлагается применить модифицированный метод динамического сгущения, а вторую задачу предполагается решить с помощью метода нахождения локального оптимума.

Образуем исходную матрицу. Исходные данные, состоящие из $\Omega' = \{P_\tau\}$, необходимые для решения задач $Z = \{Z_j\}$ могут быть представлены матрицей информационных потребностей $M = (x_\tau^j)$ элементы которой $x_\tau^j = 1$, если $P_\tau \in \Omega_j$, и $x_\tau^j = 0$ в противном случае, где Ω_j - множество ЭИМ, которое необходимо для решения Z_j .

Как правило, слабозаполненная или разреженная матрица представляет собой матрицу с большим количеством нулей. Для хранения такой матрицы при расчетах потребуются чрезмерные затраты машинного времени и объема памяти.

Разработан метод и алгоритм исключения нулевых элементов, которые позволяют значительно снизить стоимость не только хра-

нения информации, но и стоимость затрат на вычисления.

Данный метод предполагает хранить в списке только значения индексов i и j матрицы, с ненулевыми элементами. Таким образом, он позволяет преобразовать исходную матрицу M на список значений индексов ненулевых элементов матрицы M .

Используем процедуру разбиения на классы. Разобьем исходную матрицу на блоки согласно следующих требований, т.е. определяя блоки наименьшей мощности и производя разбиение на минимальное количество блоков. Пусть, в качестве примера, исходная матрица выглядит следующим образом:

	Z_3	Z_4	Z_5	Z_6	Z_7	Z_1	Z_2	Z_8	Z_9	Z_{10}
P_3	0	0	0	0	0	1	0	1	1	0
P_6	0	0	0	0	0	0	1	0	1	0
P_7	0	0	0	0	0	0	1	1	0	1
P_4	1	0	1	0	0	1	0	1	0	1
P_5	0	1	0	1	1	0	1	0	1	0
P_8	1	0	1	0	0	1	0	1	0	1
P_9	0	1	0	0	0	1	0	0	0	1
P_1	1	0	1	0	0	1	0	1	0	1
P_2	0	1	0	1	1	0	1	0	1	0
P_{10}	0	1	0	0	1	0	1	0	0	0

Разбиение $P' = (P'_1, P'_2, P'_3)$ множества Ω' и $Q' = (Q'_1, Q'_2)$ множества Z соответственно находятся по строкам и столбцам матрицы.

Тогда задача состоит в нахождении однородных пар классов (P'_k, Q'_l) , т.е., либо заполненных единицами, либо нулями, и каждая пара (P'_k, Q'_l) классов ассоциируется с банарным идеальным значением 0 или 1.

I. Таким образом, создается бинарная матрица по K строкам и L столбцам, так называемое ядро, и обозначается как $L = (a_{kl}^e)$, где $a_{kl}^e \in \{0, 1\} \quad \forall k \in \overline{1, K}$ и $\forall l \in \overline{1, L}$; a_{kl}^e есть представление пар классов (P'_k, Q'_l) .

Уточним понятие идеального значения. Когда в паре классов (P'_k, Q'_l) в матрице M_1 стоит большее количество единиц, чем количество нулей, то $a_{kl}^e = 1$. В противном случае $a_{kl}^e = 0$.

Исходя из вышеизложенного, суть модифицированного метода динамического сгущения состоит в определении элементов, которые могли бы служить представителями всех пар классов (P'_k, Q'_l) , элементов всех классов, таких чтобы сумма расстояния от всех пар классов (P'_k, Q'_l) до их представительства была минимальной.

Итак, математическая модель модифицированного метода динамического сгущения описывается следующим образом

$$\min W(P', Q', L) = \sum_{k=1}^K \sum_{l=1}^L \sum_{z_j \in P'_k, Q'_l} |x_r^j - a_{kl}^e|, \quad (I)$$

где $x_r^j = \begin{cases} 1, & \text{если } Z_j \text{ используем ЭИМ } P'_r, \\ 0 & \text{в противном случае.} \end{cases}$

$$a_{kl}^e = \begin{cases} 1, & \text{если пара } (k, l) \text{ классов заполнена единицами,} \\ 0 & \text{в противном случае.} \end{cases}$$

Из (I) можно определить ряд $(H^n, L^n): H^n = P'$ и $L^n = L^*(P', Q')$

Пусть две функции g и f соответственно для представительства и назначения.

Функция представительства предназначена для определения ядра или представительства

$$g: H^n \rightarrow L^n, \quad \text{где } L^n = g(H^n).$$

Функция назначения предназначена для отнесения элементарных информационных модулей либо задач к классу

$$f: L^{n-1} \rightarrow H^n, \quad \text{где } H^n = f(L^{n-1})$$

Поскольку из предыдущего представительства L^{n-1} определяется следующий класс ЭИМ либо задач $H^n = f(L^{n-1})$,

$$L^n = g(H^n) = L^*(H^n, Q')$$

, то ряд стационарирует (H^n, L^n) и уменьшает значение W до стационарности.

Если будем считать, что ряд стационарный с N порядка, получим:

$$W(H^0, Q', L^0) > \dots > W(H^n, Q', L^n) = W(H^{n+1}, Q', L^{n+1}) = \dots$$

В случае достижения стационарности $N=0$, т.е. с самого начала, можно симметричным образом конструировать разбиение Q'' и ядро L' , которое уменьшает критерий W . В этом случае образуется ряд (P^n, Q'', L^n) , который выполняет соотношение

$$W(P^0, Q'', L^0) \geq W(P^1, Q'', L^1) \geq \dots \geq W(P^n, Q'', L^n) \geq \dots \quad (2)$$

Алгоритм модифицированного метода конкретизируется двумя промежуточными этапами.

Можно выразить функцию $W(P', Q', L)$ меры сходства:

$$W(P', Q', L) = \sum_k \sum_{q \in Q_k} \sum_{\ell} \sum_{z_j \in Q'_\ell} |x_{z_j}^j - a_{k\ell}^j|$$

Пусть

$$A = \sum_{z_j \in Q'_\ell} |x_{z_j}^j - a_{k\ell}^j| \text{ и } q = \text{card } Q'_\ell, \text{ т.е. мощность } Q'_\ell$$

Вспользуемся переменной y_{ℓ}^j , где $y_{\ell}^j = \sum_{z_j \in Q'_\ell} x_{z_j}^j$. Это означает, что вычисляется только сумма $\sum_{z_j \in Q'_\ell} x_{z_j}^j$ всех единиц в классе. Тогда A выражается в таком виде: $A = |y_{\ell}^j - q_{\ell} a_{k\ell}^j|$

$$\text{поскольку } A = \begin{cases} y_{\ell}^j, & \text{если } a_{k\ell}^j = 0, \\ q_{\ell} - y_{\ell}^j, & \text{если } a_{k\ell}^j = 1 \end{cases}$$

В этом случае можно записать:

$$W(P', Q', L) = \sum_k \sum_{q \in Q_k} \sum_{\ell} |y_{\ell}^j - q_{\ell} a_{k\ell}^j| \rightarrow \min \quad (3)$$

Выразим через y_K^e матрицу (RQ) . Все ядра представляются в виде $(q_1 a_{K1}^e, q_2 a_{K2}^e, \dots, q_L a_{KL}^e)$, где $a_{Kl}^e \in \{0, 1\}$. Необходимо находить для всех классов элементы $(q_1 a_{K1}^e, q_2 a_{K2}^e, \dots, q_L a_{KL}^e)$ которые минимизируют

$$\sum_{P_K^e \in P_K^e} \sum_{\ell} |y_{\ell}^e - q_{\ell} a_{K\ell}^e| = \sum_{\ell} \sum_{P_K^e \in P_K^e} |y_{\ell}^e - q_{\ell} a_{K\ell}^e|. \quad (4)$$

Это приводит к нахождению элементов $q_{\ell} a_{K\ell}^e$, которые минимизируют $\sum_{P_K^e \in P_K^e} |y_{\ell}^e - q_{\ell} a_{K\ell}^e|$.

Если $a_{K\ell}^e = 0$, то мы получим $\sum_{P_K^e \in P_K^e} y_{\ell}^e$.

Если $a_{K\ell}^e = 1$, тогда

$$\sum (q_{\ell} - y_{\ell}^e) = s_K q_{\ell} - \sum y_{\ell}^e, \quad s_K = \text{card}(P_K^e), \quad \text{т.е.}$$

мощность P_K^e .

Таким образом, можно определить значения ядра. Если $a_{K\ell}^e = 0$, то сумма элементов в (P_K^e, Q_{ℓ}^e) ближе к нулю и $a_{K\ell}^e = 1$, если эта сумма приближается к $s_K q_{\ell}$.

В общем виде можно выразить A таким образом

$$A = \sum_{P_K^e \in P_K^e} \sum_{z_j \in Q_{\ell}^e} |x_{\ell}^j - a_{K\ell}^e|.$$

Пусть

$$y_K^e = \sum_{P_K^e \in P_K^e} \sum_{z_j \in Q_{\ell}^e} x_{\ell}^j, \quad \text{тогда}$$

$$A = \begin{cases} y_K^e, & \text{если } a_{K\ell}^e = 0 \\ s_K q_{\ell} - y_K^e, & \text{если } a_{K\ell}^e = 1 \end{cases} \quad \text{или} \quad A = |y_K^e - s_K q_{\ell} a_{K\ell}^e| \quad (5)$$

Тогда задача состоит в определении $s_K q_{\ell} a_{K\ell}^e$ для всех классов, которые минимизируют

$$|y_K^e - s_K q_{\ell} a_{K\ell}^e|.$$

В результате для нашего примера получим

	Z_1	Z_3	Z_7	Z_8	Z_{10}	Z_2	Z_4	Z_6	Z_7	Z_9
P_3	I			I					I	
P_4				I	I	I				
P_9	I				I		I			
P_1	I	I	I	I	I					
P_3	I	I	I	I	I					
P_4	I	I	I	I	I					
P_2						I	I	I	I	I
P_5						I	I	I	I	I
P_6						I	0	0	0	I
P_{10}						I	I	0	I	0

$$L = \begin{pmatrix} 0 & 0 \\ I & 0 \\ 0 & I \end{pmatrix}$$

Отметим, что альтернативное множество ИМ представляет собой множество ИМ, полученных с помощью некоторых операций, имеющее меньшую мощность по сравнению с первоначальным множеством возможных вариантов ИМ. Для построения альтернативного ИМ, имеется пара смежных классов задач Q'_{i-1} и Q'_i из матрицы M_2 . Пусть S_i - множество ЭИМ покрывающих информационные потребности класса задач Q'_i , элементный состав которого определяется столбцами матрицы M_2 . Каждая пара смежных классов задач порождает некоторые информационные модули. Несколько из них могут быть получены с помощью операции объединения (\cup), пересечения (\cap), разности (\setminus) и дополнения \bar{S}_i .

Используя эти операции, можно выражать одни множества через другие. Тогда пара смежных классов задач порождает $S_i, S_{i-1}, S_i \cap S_{i-1}, S_{i-1} \setminus S_i, S_i \setminus S_{i-1}, S_i \cup S_{i-1}, \bar{S}_i$ и \bar{S}_{i-1} .

Можно выражать $\bar{S}_{\ell-1}$ и \bar{S}_{ℓ} через $S_{\ell-1} \setminus S_{\ell}$ и $S_{\ell} \setminus S_{\ell-1}$. Исключая $S_{\ell-1} \cup S_{\ell}$ в связи с тем, что построение информационных массивов в виде единого набора ЭИМ приводит к практически полному отсутствию гибкости структуры баз данных.

Такая пара смежных классов задач $(Q'_{\ell-1}, Q'_{\ell})$ порождает некоторое множество ИМ, которые могут быть включены в альтернативное множество:

$$\Omega'_i = \{S_{\ell-1}; S_{\ell}; S_{\ell-1} \cap S_{\ell}; S_{\ell-1} \setminus S_{\ell}; S_{\ell} \setminus S_{\ell-1}\}.$$

Таким образом, рассмотренный подход построения множества альтернативных информационных модулей позволяет значительно уменьшить возможные варианты ИМ с 2^{τ} до $4(L-1)+1$.

Для данного тестового примера $\tau=10$ и $L=2$. Тогда уменьшение идет от 1024 вариантов ИМ до 5 вариантов. В дальнейшем выбор наилучшего варианта по составу, количеству и качеству ИМ осуществляется из 5 вариантов с помощью метода нахождения локального оптимума, что значительно сокращает затраты на время.

Введем понятие так называемого идеального информационного модуля, который содержит те и только те элементарные информационные модели, которые необходимы для Z_j .

Далее введем величину C_j для оценки ИМ F'_j , $j' \in \bar{J}'$. Отметим, что F'_j определяет j' -й ИМ.

Введем следующие величины:

$\tau'_{jj'} = P'_{jj'} / P_j$, $\tau''_{jj'} = P''_{jj'} / P_j$, P_j - мощность Z_j (необходимые данные для Z_j), $P_{jj'}$ - мощность множества ЭИМ отсутствующих в F'_j и необходимых для Z_j , $P''_{jj'}$ - мощность множества лишних для Z_j ЭИМ, входящих в состав F'_j . Тогда величина $\tau_{jj'} = \tau'_{jj'} + \tau''_{jj'}$ будет определять оценку отклонения эле-

ментного состава ИМ F_j , от идеального ИМ для Z_j . Суммарная оценка F_j' всех задач выражается в виде $C_j' = \tau_j' = \sum_{j=1}^I \tau_{jj}'$ представляет сумму всех лишних и отсутствующих нужных ЭИМ в j -ом ИМ при решении задач $\{Z_j\}$.

Тогда задача оптимизации информационных массивов, исходя из указанных выше требований к выделяемым из альтернативного множества Ω_1' , можно выразить требованием минимизации критерия качества C_j' .

Итак, математическая модель структуризации информационных массивов можно представить в виде следующей задачи линейного программирования с булевыми переменными

$$\Phi = \sum_{j=1}^{I'} c_j x_j \quad (6)$$

при ограничении $\sum_{j=1}^{I'} b_{rj} x_j = 1, \quad x = \{0, 1\}$.

Из множества Ω_2' выбираем непересекающиеся ИМ, т.е. из

$$\Omega_2' = \{S_{1-1}, S_2, S_1 \cap S_2, S_1 \setminus S_2, S_2 \setminus S_1\}, \quad \ell = 1, 2$$

Тогда получим

$$\Omega'' = \{S_1, S_2 \setminus S_1, \{S_2, S_1 \setminus S_2\}, \{S_1 \setminus S_2, S_1 \cap S_2, S_2 \setminus S_1\}\}$$

Из всех m возможных вариантов ИМ необходимо ввести выбор по k непересекающимся ($k < m$) ИМ для покрытия информационной потребности. Каждый ИМ имеет информационное качество C_j .

Пусть множество $J'' \subset \{1, 2, 3, \dots, m\} \quad |J''| = k$ - его мощность.

Тогда задача оптимизации состоит в нахождении J'' так, чтобы

$$\Phi = f(J'') = \sum_{j=1}^m c_j x_j \rightarrow \min \quad (7)$$

при ограничениях

$$\sum_{j=1}^m b_{rj} x_j = 1, \quad x = \{0, 1\}$$

Далее, имеется множество J'' , $\forall j' \in J''$. Найдем, если су-

существует, лучшую оценку, то есть, если существует $j'' \in \{1, 2, 3, \dots, m\} \setminus \mathcal{J}''$, такое, что $f(\mathcal{J}'' \cup \{j''\} \setminus \{j'\}) < f(\mathcal{J}'')$. Если такой случай проверен, то меняем \mathcal{J}'' на $\mathcal{J}_1 = \mathcal{J}'' \cup \{j''\} \setminus \{j'\}$ и процедуру повторяем до тех пор, пока найдем локальный оптимум.

В результате получим $\{S_1, S_2 \setminus S_1\}$. Тогда в состав базы данных входят ИМ $\{S_1, S_2 \setminus S_1\}$. Таким образом, метод нахождения локального оптимума позволяет найти набор ИМ, который формирует базу данных с самой минимальной суммой всех лишних и отсутствующих нужных ЭИМ при решении всех задач $\{Z_j\}$.

Для создания логической структуры баз данных кроме необходимых и достаточных используемых элементов всеми решаемыми задачами, надо знать отношения между ними. Выбор системы отношений осуществляется с помощью модифицированного метода выбора системы отношений между элементами информации.

Пусть каждый пользователь декларирует свою систему отношений U_m^j , где m - номер пользователя и j - номер задачи $/ j = \overline{1, T}, m = \overline{1, M} /$ в задаваемом им множестве элементов информации. Причем, среди этих отношений могут быть эквивалентные, независимые и пересекающиеся. Поэтому возникает необходимость анализа исходных информационных структур, представляемых пользователем, с целью их совмещения в единой структуре.

Исходные структуры пользователей можно представлять в виде ориентированных графов. Если каждому элементу ИМ S_1 и $S_2 \setminus S_1$ идентифицировать некоторую вершину графа X_i , геометрический образ которой на графе будет задан точкой, то совокупность их определенная на графе с помощью отрезков кривых,

идентифицирует U_j . Они образуют множество отношений $U_j \in U$. Совокупность элементов ИМ $S_1, S_2 \setminus S_1$ составляют элементы множества X , которое вместе с множеством U представляют собой граф $G(X, U)$, моделирующий структуру рассматриваемой системы.

Из каждой системы отношений U_m^i можно образовать множество декларируемых отношений A_m^i .

$$A_m = \bigcup_{j=1}^r A_m^j$$

Ориентированные графы $G_m^i(X_m^i, U_m^j)$ соответствуют матрице смежностей $\| \tau_{i'j'}^{A_m^i} \|$

$$\tau_{i'j'}^{A_m^i} = \begin{cases} 1, & \text{если существует ориентированная дуга} \\ & \text{в множестве отношений, соединяющая} \\ & i' \text{ с } j'; \\ 0 & \text{в противном случае.} \end{cases}$$

Эти графы могут быть связными и содержать одни и те же отношения. Поэтому преобразуем исходные графы пользователей следующим образом.

Приводим все матрицы $\| \tau_{i'j'}^{A_m^i} \|$ к одному размеру и формируем матрицу $\| R_{i'j'}^j \|$ того же размера.

$$R_{i'j'}^j = \sum_{m=1}^M \tau_{i'j'}^{A_m^i}, \quad \text{где } m = \overline{1, M}; j = \overline{1, J}.$$

Для каждой матрицы $\| R_{i'j'}^j \|$ формируем матрицу $\| \tau_{i'j'}^j \|$ таким образом, чтобы

$$\tau_{i'j'}^j = \begin{cases} 1, & \text{если } R_{i'j'}^j \geq 1, \\ 0 & \text{в противном случае.} \end{cases}$$

Устраняем транзитивные отношения таким образом. Определим все транзитивные отношения: там, где стоят единицы, ставим нули. Так например, имеются три элемента x_1, x_2, x_3 . Они формируют такие отношения как $(x_1, x_2), (x_2, x_3), (x_1, x_3)$. Из этих отношений (x_1, x_3) является лишним. Тогда в матрице смежности ставим нуль для (x_1, x_3) и получим матрицу $\|r_{ij}^j\|$. Матрица сложности окончательной синтезированной логической структуры БД определяется

$$R_{ij'} = \sum_{j=1}^n r_{ij'}^j.$$

Формируем для матрицы $R_{ij'}$ также матрицу $\|R'_{ij'}\|$ таким образом, чтобы

$$R'_{ij'} = \begin{cases} 1, & \text{если } R_{ij'} \geq 1 \\ 0 & \text{в противном случае.} \end{cases}$$

Для всех задач $Z_j, j = \overline{1, T}$, условие независимости описывается следующим образом.

$$\text{card}((A_3 \cup \dots \cup A_m) \cap (A_k \cup \dots \cup A_T)) \leq \Delta_n,$$

где $A_1, A_m, \dots, A_k, \dots, A_T$ - множества декларируемых отношений для всех решаемых задач.

Δ_n - некоторая мера независимости множеств отношений, которая может быть задана, исходя из необходимой точности последующих вычислений.

В результате применения таково модифицированного метода выбора системы отношений получим все отношения, приведенные в табл. I.

Третья глава посвящена выделению классов однотипных элементов, реализации элементарных запросов, разработке метода и алгоритма синтеза оптимальной структуры баз данных с помощью метода исключения бесперспективных вариантов, а также постро-

Таблица I.

Необходимые и достаточные отношения для создания
логической структуры БД фрагментов предметных областей
о заработной плате рабочих.

№ отно- шения	Тип отно- шения	Начало отношения	Конец отношения	Отношения
1	I:I	TN	F ₁₀	(TN, F ₁₀)
2	I:I	TN	K	(TN, K)
3	I:I	TN	C	(TN, C)
4	I:I	IZ	NIZ	(IZ, NIZ)
5	I:I	IZ	KZ	(IZ, KZ)
6	I:I	IZ	NV	(IZ, NV)
7	I:M	SM	U	(SM, U)
8	I:M	P _Д	B _Д	(P _Д , B _Д)
9	I:I	F ₁₀	B _Д	(F ₁₀ , B _Д)
10	I:I	SM	KZ	(SM, KZ)
11	I:I	SM	NV	(SM, NV)
12	I:I	IZ	KOP	(IZ, KOP)
13	I:M	PR	K	(PR, K)
14	I:I	K	WO	(K, WO)
15	I:I	SI	RC	(SI, RC)
16	I:I	IZ	KZAT	(IZ, KZAT)
17	I:I	NIZ	SI	(NIZ, SI)
18	I:M	C	U	(C, U)
19	I:I	C	NC	(C, NC)
20	I:I	TN	SM	(TN, SM)
21	I:I	IZ	KP	(IZ, KP)
22	I:I	KRZ	K	(KRZ, K)
23	I:I	K	P _Д	(K, P _Д)

ению оптимальной логической структуры. В главе показано, что целесообразно выделить классы однотипных элементов, таких, что структуры отношений элементов принадлежат к одному классу. Для этого используется метод динамического сгущения.

Имеется множество Ω''' , в состав которого входят элементы базы данных. образуем два непересекающихся подмножества множества Ω''' , \mathcal{D} и V , мощностью которых являются n атрибутов и p количественных переменных соответственно. Если предположить, что все элементы находятся в матрице X , тогда элементы множества \mathcal{D} расположены по n строкам и элементы множества V по p столбцам.

$$X = (x_{ij}^j), \quad i \in \mathcal{D} \quad \text{и} \quad j \in V,$$

где x_{ij}^j - значение переменной j для атрибута i .

Предполагается, что каждому атрибуту приписан некоторый вес p_i , сумма которых равна $\sum_{i \in \mathcal{D}} p_i = 1$, где

$$p_i = \frac{1}{n} \quad \forall i \in \mathcal{D}.$$

Предположим также, что каждый атрибут i можно представить как точку. Тогда, для атрибута i можно ассоциировать вектор $x_i = (x_{i1}^1, \dots, x_{ip}^p) \in R^p$. Множество векторов x_i с весами p_i образуют облако $N(\mathcal{D}) \in R^p$.

Предположим, что матрица X , центрированная по столбцам, то есть, сумма элементов каждого столбца равна нулю. В противном случае необходимо ввести преобразование для выполнения этого условия. Это требование получается в результате того, что центр тяжести облака $N(\mathcal{D})$ находится в начале пространства R^p .

$$\sum_{i \in \mathcal{D}} p_i x_i^j = 0, \quad \forall j \in V.$$

Также для каждой переменной j аналогичным образом ассоциируем ее с вектором $x^j = (x_1^j, \dots, x_n^j)$. Множество x^j с весами q_j образуют облако $N(V) \in R^n$.

Имеется пространство атрибутов с квадратичной метрикой, определенной диагональной матрицей, где коэффициенты важности q_j находятся по диагональной линии. Тогда

$$d^2(x_i, x_{i'}) = \sum_{j \in V} q_j (x_i^j - x_{i'}^j)^2.$$

В этом выражении могут существовать значения переменной j с разными единицами измерения. Чтобы каждая переменная j сыграла одинаковую роль, необходимо ввести нормализацию коэффициентов важности q_j таким образом:

$$q_j = \frac{1}{\sigma_j^2}, \quad \text{где } \sigma_j^2 = \sum_{i \in B} p_i (x_i^j)^2.$$

σ_j - представляет собой дисперсию переменной j .

Имеется пространство переменных с весами \mathcal{D}_p и метрика. Тогда

$$d^2(x^j, x^{j'}) = \sum_{i \in B} p_i (x_i^j - x_i^{j'})^2.$$

Введем понятие инерции. Рассмотрим облако из n точек (x_1, x_2, \dots, x_n) пространства R^p , причем каждой точке x_k приписан некоторый вес μ_k ($k = 1, 2, \dots, n$). Инерцией множества (x_1, x_2, \dots, x_n) вокруг точки $a \in R^p$ называется величина

$$I(a) = \sum_{k=1}^n \mu_k d^2(x_k, a).$$

Вклад каждой точки в эту инерцию есть $\mu_k d^2(x_k, a)$.

$$I(G) = \sum_{k=1}^n \mu_k d^2(x_k, G) = \sum_{j=1}^{K'} \sum_{x_k \in C_j} \mu_k d^2(x_k, G_j) + \sum_{j=1}^{K'} m_j d^2(G_j, G), \text{ где } G_j = \sum_{x_k \in C_j} \mu_k x_k / \sum_{x_k \in C_j} \mu_k -$$

центры тяжести K' классов. Эти центры образуют облака из K' новых объектов, которым соответствуют веса

$$\sum_{j=1}^{K'} m_j G_j / \sum_{j=1}^{K'} m_j \quad m_j = \sum_{x_k \in C_j} \mu_k \quad - \text{ центр тяжести облака } \{x_k\}, k = \overline{1, n}$$

с весами $\{\mu_k\}$, $k = \overline{1, n}$ является также центром тяжести облака $\{G_j\}$, имеющегося веса $\{m_j\}$, $j = \overline{1, K'}$.

Инерция облака $I(G)$ вокруг центра тяжести также называется полной инерцией облака и обозначается T . Величина $\sum_{j=1}^{K'} m_j d^2(G_j, G)$ есть инерция облака центров тяжести $\{G_j\}$. Ее называют межклассовой инерцией и обозначают B . Каждое выражение $\sum_{x_k \in C_j} \mu_k d^2(x_k, G_j)$ представляет инерцию точек, составляющих класс C_j , вокруг центра тяжести этого класса. Сумма этих K' инерции называется внутриклассовой инерцией и обозначается W .

Уравнение разложения инерции записывается как

$$T = W + B.$$

Отметим, что метод приводит к минимизации внутриклассовой инерции и наоборот, к максимизации межклассовой инерции.

Определяется инерция облака $N(Z)$ по отношению к началу координаты, которая представляет также центр тяжести:

$$I(Z, p_i, q_i) = \sum_{i \in Z} p_i d^2(0, x_i) = \sum_{i \in Z} p_i \sum_{j \in V} q_j (x_i^j)^2 = \sum_{i \in Z} \sum_{j \in V} p_i q_j (x_i^j)^2.$$

Далее, из исходной матрицы X сформируем новую матрицу, в которой по строкам расположено разбиение множества \mathcal{D} на непустые подмножества $P' = (P'_1, \dots, P'_{K'})$ на K' классов и по столбцам разбиения множества V на непустые подмножества $Q' = (Q'_1, \dots, Q'_L)$ на L классов. Элементы матрицы $X(P', Q')$ определяются таким образом:

$$x_{kl}^e(P', Q') = \sum_{i \in P'_k} \sum_{j \in Q'_l} p_i q_j x_i^j / \sum_{i \in P'_k} \sum_{j \in Q'_l} p_i q_j.$$

С новой матрицей $X(P', Q')$ связываются некоторые веса $\gamma_{k'} = \sum_{i \in P'_k} p_i$, $\forall k' = 1, \dots, K'$ и $\beta_l = \sum_{j \in Q'_l} q_j$, $\forall l = 1, \dots, L$.

Исходя из ранее изложенного, задача нахождения однотипных элементов по классам в системе состоит в нахождении разбиения P' и Q' множеств \mathcal{D} и V соответственно такого, чтобы новая матрица $X(P', Q')$, имеющая веса $(\gamma_{k'})$ и (β_l) , потеряла меньше информации, то есть

$$I(X(P', Q'), (\gamma_{k'}), (\beta_l)) = \sum_{k'} \sum_l \gamma_{k'} \beta_l (x_{k'l}^e)^2 \rightarrow \max.$$

Таблица 2

Расчет количественных атрибутов

$i \backslash j$	KZ	KP	NV	RC	KOLD	KZAT	SI
I	2	3	4	5	6	7	8
TN	15	8	5	0	0	0	0
F ₁₀	15	0	0	0	0	0	0
K	0	0	0	0	0	0	0
PR	0	0	0	0	0	0	0
KRZ	0	0	0	0	0	0	0
C	1500	1000	40	0	0	0	0

I	2	3	4	5	6	7	8
U	150	92	40	0	0	0	0
SM	150	92	5	0	0	0	0
WO	0	0	0	0	0	0	0
NC	1500	100	40	0	0	0	0
IZ	150000	100000	90	3200	30	1500	15
NIZ	150000	100000	30	3200	30	1500	15
KOP	0	0	0	0	0	0	0
PD	0	0	0	5	0	4	0
BD	0	0	0	45	0	41	0
Сумма	303330	201292	250	6450	60	3045	30
\bar{X}	20222	13419,41	16,67	430	4	203	2

Таблица 3.

Результативная таблица.

$i \backslash j$	KZAT	KOLD	SI	RC	NV	KP	KZ
I	2	3	4	5	6	7	8
TN	0,40	0,39	0,39	0,40	0,46	0,39	0,40
F10	0,40	0,39	0,39	0,40	0,66	0,40	0,40
SM	0,40	0,39	0,39	0,40	0,46	0,39	0,39
C	0,40	0,39	0,39	0,40	-0,92	0,37	0,37
U	0,40	0,39	0,39	0,40	-0,92	0,39	0,39
NC	0,40	0,39	0,39	0,40	-0,92	0,39	0,37
KRZ	0,40	0,39	0,39	0,40	0,66	0,40	0,40

I	2	3	4	5	6	7	8
IZ	-2,55	-2,55	-2,55	-2,55	-2,89	-2,55	-2,55
NIZ	-2,55	-2,55	-2,55	-2,55	-0,52	-2,55	-2,55
KOP	0,40	0,39	0,39	0,40	0,66	0,40	0,40
PR	0,40	0,39	0,39	0,40	0,66	0,40	0,40
K	0,40	0,39	0,39	0,40	0,66	0,40	0,40
WO	0,40	0,39	0,39	0,40	0,66	0,40	0,40
BD	0,32	0,39	0,39	0,35	0,66	0,40	0,40

$$I = 0,0000305; \quad q_1 = 0,0000196; \quad q_2 = 0,0000294;$$

$$q_3 = 0,393580; \quad q_4 = 0,0009204; \quad q_5 = 0,0980581;$$

$$q_6 = 0,0019653; \quad q_7 = 0,1961161.$$

Тогда все классы однотипных элементов описываются таким образом:

$$P'_1 = \{TN, F10, SM\}, \quad P'_2 = \{C, U, NC, KRZ\}, \quad P'_3 = \{IZ, NIZ, KOP\}.$$

$$P'_4 = \{PR, K, WO, PD, BD\}; \quad Q'_1 = \{KZAT, KOLB, SI\}.$$

$$Q'_2 = \{RC, NC, KP, KZ\}.$$

Реализация элементарных запросов осуществляется следующим образом.

Пусть имеются некоторые множества отношений $\{A_j\}$ $j = \overline{1, J}$, необходимые для решения задач $\{Z_j\}$.

Тогда проблема состоит в формировании связанного подграфа для каждого множества A_j , чтобы для любого его элемента существовали пути /дуги/ от него к остальным элементам этого множества.

В этой главе разработан метод и алгоритм синтеза оптимальной логической структуры базы данных.

Логическая структура представляет собой множество элементов и отношений между ними. Тогда можно представить ее в виде ориентированного графа $G = (X, U)$, где $X = \{x_i\}$, $i = \overline{1, n}$ - множество вершин, $U = \{\mu_j\}$, $j = \overline{1, m}$ - множество дуг, соединяющих между собой пары вершины. Каждая вершина x_i графа G однозначно сопоставляется информационными элементами определенного типа, каждая дуга μ_j - связь между двумя какими-либо информационными элементами.

Итак, модель синтеза оптимальной логической структуры БД описывается следующими соотношениями.

Определить множество дуг $\{x_{IT}^{ijk} = 1\}$, входящих в оптимальное решение при минимальной годовой стоимости поиска, хранения и передачи информации

$$W = \sum_{k=1}^K \sum_{j=1}^P \sum_{i=1}^P \sum_{l=1}^n \sum_{j=1}^n a f_{ijk} x_{IT}^{ijk} \rightarrow \min,$$

где a - средняя стоимость обработки одного указателя связи /включает стоимость преобразования ключа в адрес, стоимость поиска, чтения соответствующего элемента информации и передачи информации; f_{ijk} - частота связи (i, j) в K -ом запросе; x_{IT}^{ijk} - бинарная переменная, используемая запросами.

$$x_{IT}^{ijk} = \begin{cases} 1, & \text{если дуга } (i, j) \in I \times J \text{ в } K\text{-ом запросе;} \\ 0, & \text{в противном случае.} \end{cases}$$

Принимая во внимание следующие условия

$$\lim_{n \rightarrow +\infty} \left(S + \frac{\tau}{2n} \right) \rightarrow S,$$

где S - минимальная стоимость хранения логической структуры базы данных;

r - радиус окружности;

n - искомый параметр.

$i \leq j$ и $i < j$ поскольку граф ориентированный и без петель.

Приведенная модель синтеза оптимальной логической структуры БД легко сводится к задаче целочисленного линейного программирования.

Пронумерованы классы и их элементы, полученные в результате действия метода динамического сгущения

$$1 = \{1, 2, 3\}, \quad 2 = \{4, 5, 6, 7\}, \quad 3 = \{8, 9, 10\},$$

$$4 = \{11, 12, 13, 14, 15\}, \quad 5 = \{16, 17, 18\}, \quad 6 = \{19, 20, 21, 22\},$$

где $TN = 1$; $FO = 2$; $SM = 3$; $C = 4$; $V = 5$; $NC = 6$;

$KRZ = 7$; $IZ = 8$; $NIZ = 9$; $KOP = 10$; $K = 11$; $PR = 12$;

$WO = 13$; $PD = 14$; $BD = 15$; $KZAT = 16$; $KOLD = 17$;

$SI = 18$; $RC = 19$; $NV = 20$; $PK = KP = 21$; $KZ = 22$.

Таблица 4.

Частота решения задач.

Z_j	f_j	Z_j	f_j
1	180	6	30
2	150	7	20
3	100	8	10
4	80	9	10
5	50	10	5

В результате действия алгоритма синтеза оптимальной логической структуры базы данных получается для нашего примера логическая структура со стоимостью хранения $U = 960$ руб. и суммарной стоимостью поиска и передачи информации $W = 10760$ руб.

Элементарные запросы реализуются таким образом.

- { /8,9/, /8,20/, /8,22/, /3,22/, /1,3/, /1,4/, /1,5/, /1,2/, /1,11/, /1,12/ } ;
- { /3,20/, /3,22/, /3,5/, /1,3/, /1,2/, /2,15/, /1,15/, /14,15/ } ;
- { /8,9/, /8,17/, /8,22/, /3,22/, /1,3/, /8,10/, /1,2/, /10,17/ } ;
- { /8,9/, /9,18/, /18,19/, /8,20/, /3,20/, /1,3/, /1,11/, /7,11/, /11,14/, /11,12/, /11,13/ } ;
- { /8,9/, /8,16/, /9,18/, /8,19/, /8,18/ } ;
- { /8,9/, /8,20/, /8,21/, /3,21/, /1,3/, /1,2/, /4,5/, /4,6/, /1,4/, /1,5/ } ;
- { /8,9/, /4,5/, /8,16/, /8,19/, /8,18/, /3,20/, /3,5/ } ;
- { /8,9/, /8,22/, /1,3/, /1,2/, /1,11/, /1,12/, /1,13/ } ;
- { /8,9/, /8,16/, /8,20/, /1,3/, /1,2/ } ;
- { /1,4/, /4,6/ } .

В этой главе также построена оптимальная логическая структура. Она представлена в виде гиперграфов на рис. 1.

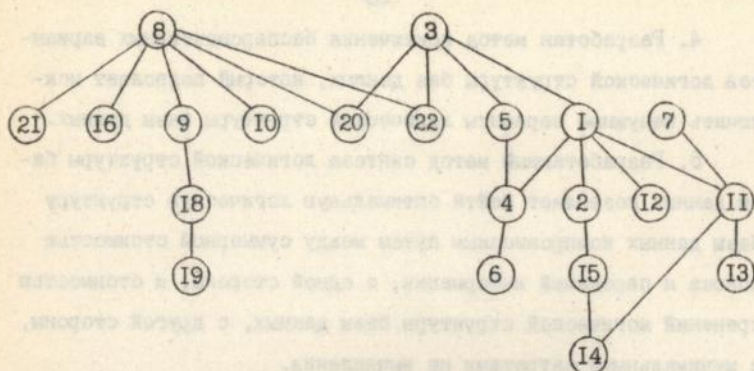


Рис. 1. Логическая структура в виде гиперграфа.

Основные выводы.

1. Анализ оптимизации распределения информационных потоков с применением метода сортировки, показывает, что этот метод не позволяет устранить дублирование информации, что приводит к неоптимальной логической структуре баз данных. Для устранения этого недостатка применен модифицированный метод динамического сгущения.

2. Исследование метода динамического сгущения показывает, что он не только не позволяет устранить дублирование информации в покрытии, но и не может найти предельное количество попарных пересечений классов, нужных для нахождения разбиения.

Модифицированный метод динамического сгущения позволяет устранить эти недостатки.

3. Разработан также метод исключения нулевых элементов, который позволяет снизить стоимость хранения и вычисления.

4. Разработан метод исключения бесперспективных вариантов логической структуры баз данных, который позволяет исключить ненужные варианты логической структуры базы данных.

5. Разработанный метод синтеза логической структуры базы данных позволяет найти оптимальную логическую структуру базы данных компромиссным путем между суммарной стоимостью поиска и передачей информации, с одной стороны, и стоимостью хранений логической структуры базы данных, с другой стороны, с минимальными затратами на вычисления.

6. Разработанный модифицированный метод выбора системы отношения учитывает наличие всех решаемых задач, а также позволяет найти логическую структуру базы данных в виде гиперграфа. Такая новая модель базы данных использует ассоциативный метод доступа.

Основные опубликованные работы, отражающие содержание диссертации

1. Волколупова Р.Т., Амана Н.И. Идентификация подструктур баз данных. Сборник научных трудов "Экономические проблемы работы предприятий в новых условиях хозяйствования". - К.: УМКВО, 1988. - С.52-61.

2. Волколупова Р.Т., Амана Н.И. Модифицированный метод динамического сгущения. Материалы докладов учредительной конференции международной ассоциации по нетрадиционным методам оптимизации. - Красноярск: КИКТ, 1992.

3. Волколупова Р.Т., Бузько Я.П., Амана Н.И. Метод нахождения локального оптимума. Сборник научных трудов "Использование математических методов и информационных технологий в технических и экономических системах". - К.: ИК АН Украины, 1992.

Подп. к печ. 20.05.93. Формат 60×84^{1/16}.
Бумага тип. № 3 . Способ печати офсетный. Услови. печ. л. 1,63
Услови. кр.-отг. 1,86 . Уч.-изд. л. 1,0
Тираж 100 . Зак. № 4544 . Бесплатно.

Фирма «ВИПОЛ»
252151, г. Киев, ул. Вольнская, 60.

08

AB 27.603