

Не функциональные зависимости

СЛУЖ ПРО ЗАДАН

СЛУ N	ПРО N	СЛУ ЗАДАН
34	1	A
34	1	B
34	2	A
34	2	B
41	1	A
41	1	D

- служащий может участвовать в нескольких проектах
- в каждом проекте служащий выполняет одинаковый набор заданий

Ключ отношения {СЛУ_N, ПРО_N, СЛУ_ЗАДАН}

Находится в BCBF (отсутствуют нетривиальные FD)

$IF(\langle c, p1, z1 \rangle \in SP3 \text{ AND } \langle c, p2, z2 \rangle \in SP3) THEN(\langle c, p1, z2 \rangle \in SP3 \text{ AND } \langle c, p2, z1 \rangle \in SP3)$

Аномалии обновления:

- **Добавление кортежа.** Если служащий присоединяется к проекту, то необходимо добавить столько кортежей, сколько заданий выполняет этот служащий.
- **Удаление кортежей.** Если служащий прекращает участие в проектах, то отсутствует возможность сохранить данные о заданиях, которые он может выполнять.
- **Модификация кортежей.** При изменении одного из заданий служащего необходимо изменить значение атрибута СЛУ_ЗАДАН в нескольких кортежах, для каждого проекта.

Возможная декомпозиция:

СЛУЖ ПРО

СЛУ N	ПРО N
34	1
34	2
41	1

СЛУЖ ЗАДАН

СЛУ N	СЛУ ЗАДАН
34	A
34	B
41	A
41	D

В отношении нет FD, однако:

1. Декомпозиция без потерь
2. Проблемы с аномалиями решены

Многозначная зависимость.

Впервые описаны Ронем Фейджином в 1971 г. Фейджин назвал зависимости этого вида многозначными (*multi-valued dependency* – MVD).

В СЛУЖ_ПРО_ЗАДАН выполняются две MVD: СЛУ_N →→ ПРО_N и СЛУ_N →→ СЛУ_ЗАДАН

В терминах алгебраических выражений MVD СЛУ_N →→ ПРО_N означает, что результат:
 $R_i = (\text{СЛУЖ_ПРО WHERE (СЛУ_N = 'N' AND СЛУ_ЗАДАН = 'Z_i')}) \text{PROJECT \{ПРО_N\}}$
'N' = const & $\forall Z_i, Z_j \Rightarrow R_i = R_j$

В переменной отношения r с атрибутами A, B, C (в общем случае, составными) имеется **многозначная зависимость** B от A ($A \rightarrow\rightarrow B$) в том и только в том случае, когда множество значений атрибута B, соответствующее паре значений атрибутов A и C, зависит от значения A и не зависит от значения C.

IF (<a, b, c'> AND <a, b', c>) THEN (<a, b, c>)

В силу симметрии атрибутов мы можем также записать правило в виде:

IF (<a, b, c'> AND <a, b', c>) THEN (<a, b, c> AND <a, b', c'>)

Лемма Фейджина о «двойственности» многозначной зависимости:

В отношении $r \{A, B, C\}$ выполняется $MVD A \twoheadrightarrow B$ в том и только в том случае, когда выполняется $MVD A \twoheadrightarrow C$.

Доказательство достаточности условия леммы.

Пусть $\exists MVD A \twoheadrightarrow B$.

r - отношение,

$\forall a : a (A \in Hr) \in t$ - значение атрибута $A \in Hr$ в некотором кортеже тела r ,

$\{b\}$ – множество значений атрибута $B \in Hr$, взятых из всех кортежей тела Vr , в которых значением атрибута A является a .

Предположим, что для этого значения a $MVD A \twoheadrightarrow C$ не выполняется.

Это означает, что существуют такое допустимое значение c атрибута C и такое значение $b \in \{b\}$, что кортеж $\langle a, b, c \rangle \notin Vr$.

Но это противоречит наличию $MVD A \twoheadrightarrow B$. Следовательно, если выполняется $MVD A \twoheadrightarrow B$, то выполняется и $MVD A \twoheadrightarrow C$.

Необходимость условия леммы доказывается симметрично.

$MVD A \twoheadrightarrow B$ и $A \twoheadrightarrow C$ представляют в форме $A \twoheadrightarrow B|C$

Многозначные и функциональные зависимости

FD - частный случай MVD, когда множество значений зависимого атрибута состоит из одного элемента.

если \exists FD $A \rightarrow B$ ($r, \text{Hr}\{A, B, C\}$), то $\Rightarrow \exists$ MVD $A \twoheadrightarrow B|C$

Если $\langle a, b, c_1 \rangle \in r$ AND $\langle a, b_1, c \rangle \in r$, тогда в силу $A \rightarrow B \Rightarrow$

$b_1 = b$, то есть $\langle a, b_1, c \rangle = \langle a, b, c \rangle \Rightarrow$

$\langle a, b, c \rangle \in r \Rightarrow$

$\exists A \twoheadrightarrow B|C$ по определению.

Обратное – не верно.

Многозначная зависимость является обобщением понятия функциональной зависимости.

Теорема Фейджина о декомпозиции отношения с многозначной зависимостью

r отношение на $Hr \{A, B, C\}$

Отношение r декомпозируется без потерь на проекции $\{A, B\}$ и $\{A, C\}$ тогда и только тогда, когда для него выполняется $MVD A \twoheadrightarrow B|C$.

$r = r_1 \text{ NATURALJOIN } r_2, r_1=r \text{ PROJECT}\{AB\}, r_2=r \text{ PROJECT}\{AC\} \Leftrightarrow MVD A \twoheadrightarrow B|C$

[согласно определению $MVD A \twoheadrightarrow B : IF \langle ab1c \rangle \in r \text{ AND } \langle abc1 \rangle \in r \text{ THEN } \langle abc \rangle \in r$]

Доказательство **необходимости** условия теоремы (т.е. \Rightarrow)

Левая часть выполняется, т.е. $r = r_1 \text{ NATURALJOIN } r_2$ (декомпозиция без потерь)

Докажем что в таком случае выполняется $MVD A \twoheadrightarrow B|C$.

Пусть $\langle abc1 \rangle \in r$ и $\langle ab1c \rangle \in r \Rightarrow$

По определению проекций $\langle ab \rangle \in r_1, \langle ac \rangle \in r_2 \Rightarrow$

$\langle abc \rangle \in r_1 \text{ NATURALJOIN } r_2$, если исходная декомпозиция – без потерь \Rightarrow

$\langle abc \rangle \in r$ – а это и есть определение $MVD A \twoheadrightarrow B$.

Чтд, необходимость доказана.

Доказательство **достаточности** условия теоремы (т.е. \Leftarrow).

Пусть имеется многозначная зависимость $A \twoheadrightarrow B|C$.

Докажем, что декомпозиция $r = r_1 \text{ NATURALJOIN } r_2$ является декомпозицией без потерь.

Доказательство аналогично теореме Хитта и разбито на две части.

Первая – включение $r \subseteq r_1 \text{ NATURALJOIN } r_2$ – то есть надо доказать что любой кортеж исходного отношения принадлежит результату - доказывается так-же как у Хитта на основе определения операций проекции и естественного соединения и выполняется всегда для любого состояния r .

Докажем вторую половину, а именно что $r_1 \text{ NATURALJOIN } r_2 \subseteq r$.

Пусть $\langle abc \rangle \in r_1 \text{ NATURAL JOIN } r_2 \Rightarrow \langle ab \rangle \in r_1, \langle ac \rangle \in r_2 \Rightarrow \langle abc1 \rangle \in r$ и $\langle ab1c \rangle \in r \Rightarrow$

по определению многозначной зависимости $\Rightarrow \langle abc \rangle \in r$.

Включение доказано. **Достаточность доказана.** Теорема Фейджина полностью доказана.

Четвертая нормальная форма (4NF)

Переменная отношения r находится в **четвертой нормальной форме (4NF)** в том и только в том случае, когда она находится в BCNF, и все MVD r являются FD с детерминантами – возможными ключами отношения r .

4NF является BCNF, в которой многозначные зависимости вырождаются в функциональные.

СЛУЖ_ПРО_ЗАДАН не находится в 4NF, поскольку MVD $СЛУ_N \twoheadrightarrow ПРО_N \mid СЛУ_ЗАДАН$ не являются функциональными и детерминант не является возможным ключом отношения.

СЛУЖ_ПРО_Н и СЛУЖ_ЗАДАНИЕ находятся в BCNF поскольку не содержат MVD, отличных от FD с детерминантом – возможным ключом (нет нетривиальных FD). Поэтому они находятся в 4NF.

Многозначная зависимость MVD $A \twoheadrightarrow B \mid C$ называется **нетривиальной многозначной зависимостью**, если не существует функциональных зависимостей FD $A \rightarrow B$ или $A \rightarrow C$.

$СЛУ_N \twoheadrightarrow ПРО_N \mid СЛУ_ЗАДАН$ – пример нетривиальной многозначной зависимостью.

MVD $A \twoheadrightarrow B$ называется **тривиальной** если либо $B \in A$, либо $A \cup B = H_r$.

Тривиальная MVD всегда удовлетворяется.

При $B \in A$ она вырождается в тривиальную FD.

При $A \cup B = H_r$ MVD верна по определению.

N-декомпозируемые отношения

Будем называть **N-декомпозируемым** отношением отношение, которое может быть декомпозировано без потерь на N проекций.

СЛУЖ ПРО ЗАДАН (ключ Нг и отсутствуют нетривиальные MVD)

СЛУ N	ПРО_N	СЛУ ЗАДАН
34	1	A
34	1	B
34	2	A
41	1	A

СЛУЖ ПРО

СЛУ N	ПРО_N
34	1
34	2
41	1

ПРО ЗАДАН

ПРО_N	СЛУ ЗАДАН
1	A
1	B
2	A

СЛУЖ ЗАДАН

СЛУ N	СЛУ ЗАДАН
34	A
34	B
41	A

СЛУЖ ПРО NATURAL JOIN ПРО ЗАДАН

СЛУ N	ПРО_N	СЛУ ЗАДАН
34	1	A
34	1	B
34	2	A
41	1	A
41	1	B

$r \neq r1 \text{ NATURALJOIN } r2$

$r \neq r2 \text{ NATURALJOIN } r3$

$r \neq r1 \text{ NATURALJOIN } r3$

=> по Фейджину – нетривиальные MVD отсутствуют

СЛУЖ ПРО NATURAL JOIN СЛУЖ ЗАДАН

<34,2,B> - лишний

ПРО ЗАДАН NATURAL JOIN СЛУЖ ЗАДАН

<41,2,A> - лишний

НО

$r = r1 \text{ NATURALJOIN } r2 \text{ NATURALJOIN } r3$

Это говорит о том, что между атрибутами этого отношения имеется некоторая зависимость, но эта зависимость не является ни функциональной, ни многозначной зависимостью.

Зависимость проекции/соединения

Обозначим СПЗ = СЛУЖ_ПРО_ЗАДАН, СП = СЛУЖ_ПРО, ПЗ = ПРО_ЗАДАН, СЗ = СЛУЖ_ЗАДАН

СПЗ = СП NATURAL JOIN ПЗ NATURAL JOIN СЗ

⇔

IF (<сп> ∈ СП AND <пз> ∈ ПЗ AND <сз> ∈ СЗ) THEN <спз> ∈ СПЗ

Ограничение, обеспечивающее возможность восстановления без потерь, определенное на исходном отношении будет выглядеть следующим образом:

IF (<спз'> ∈ СПЗ AND <сп'з> ∈ СПЗ AND <с'пз> ∈ СПЗ) THEN <спз> ∈ СПЗ

Если служащий «С» участвует в проекте «П» выполняя какое-либо задание, и в проекте «П» выполняется задание «З» каким-либо служащим, и служащий «С» выполняет задание «З» в каком-либо проекте, ТО служащий «С» выполняет задание «З» в проекте «П».

Зависимость проекции/соединения (Project-Join Dependency – PJD)

Отношение r на $Hr\{A, B, \dots, Z\}$ (составные, перекрывающиеся атрибуты).

В переменной отношения r удовлетворяется **зависимость проекции/соединения** $*(AB \dots Z)$ тогда и только тогда, когда любое допустимое значение r можно получить путем естественного соединения проекций этого значения на атрибуты $AB \dots Z$.

$r = r\{A\} \text{ NATURALJOIN } r\{B\} \text{ NATURALJOIN } \dots \text{ NATURALJOIN } r\{Z\}$

При этом данная зависимость должна выполняться для любого значения тела отношения, удовлетворяющего всем наложенным на него ограничениям предметной области.

Аномалии, вызываемые наличием зависимости проекции/соединения

В СПЗ 2 \exists PJD *({СЛУ N, ПРО N}, {ПРО N, СЛУ ЗАДАН}, {СЛУ N, СЛУ ЗАДАН})

СЛУ N	ПРО N	СЛУ ЗАДАН
34	1	В
34	2	А

IF (<слз'> \in СПЗ AND <сл'з> \in СПЗ AND <с'пз> \in СПЗ) THEN <слз> \in СПЗ

- **Добавление кортежей.**

СПЗ_2 добавить <41,1,А>,

СЛУ N	ПРО N	СЛУ ЗАДАН
34	1	В
34	2	А
41	1	А

НЕ УДОВЛЕТВОРЯЕТ ограничению!

34	1	А
----	---	---

Кортежи <34,1,В>, <34,2,А> и <41,1,А> дают нам пары 34-1, 34-А, 1-А =>

Ограничение целостности требует включения в отношение картежа <34,1,А>

[добавление <34,1,А> не нарушает ограничения и не требует других добавлений]

- **Удаление кортежа.**

Если из СПЗ_2 удаляется кортеж <34,1,А>, то нужно удалить и кортеж <41,1,А>.

Поскольку в соответствии с ограничением целостности наличие второго кортежа означает наличие первого.

[удаление <41,1,А> не нарушает ограничения и не требует других удалений]

СЛУ N	ПРО_N	СЛУ ЗАДАН
34	1	В
34	2	А

СЛУЖ ПРО

СЛУ N	ПРО_N
34	1
34	2

+ <41,1,А>

СЛУЖ ПРО

СЛУ N	ПРО_N
34	1
34	2
41	1

ПРО ЗАДАН

ПРО_N	СЛУ ЗАДАН
1	В
2	А

ПРО ЗАДАН

ПРО_N	СЛУ ЗАДАН
1	В
2	А
1	А

СЛУЖ ЗАДАН

СЛУ N	СЛУ ЗАДАН
34	А
34	В

СЛУЖ ЗАДАН

СЛУ N	СЛУ ЗАДАН
34	А
34	В
41	А

СЛУЖ ПРО ЗАДАН = СЛУЖ ПРО NATURAL JOIN ПРО ЗАДАН NATURAL JOIN СЛУЖ ЗАДАН

СЛУ N	ПРО_N	СЛУ ЗАДАН
34	1	А
34	1	В
34	2	А
41	1	А

Связь многозначной зависимости и зависимости проекции/соединения

Зависимость проекции/соединения фактически является обобщением понятия многозначной зависимости, введенной нами для четвертой нормальной формы.

Теорема Фейджина (формулировка для зависимости проекции/соединения).

Отношение $r(ABC)$ удовлетворяет зависимости проекции/соединения $*(AB, AC)$ тогда и только тогда, когда имеется многозначная зависимость $A \twoheadrightarrow B|C$.

$$r(ABC) \models (r \text{ PROJECT}(AB)) \text{ NATURAL JOIN } (r \text{ PROJECT}(AC)) \Leftrightarrow *(AB, AC) \Leftrightarrow A \twoheadrightarrow B|C$$

Многозначная зависимость - частный случай зависимости проекции-соединения, т.е., если в отношении имеется многозначная зависимость, то имеется и зависимость проекции-соединения на 2 соответствующих атрибута.

Обратное (для любого количества атрибутов), неверно.

Тривиальные и нетривиальные PJD

Зависимость проекции/соединения $*(AB...Z)$ в отношении r называется **подразумеваемой возможными ключами** в том и только в том случае, когда каждый составной атрибут $AB...Z$ является суперключом r , т.е. включает хотя бы один возможный ключ r .

Зависимость проекции/соединения $*(AB...Z)$ в отношении r называется **тривиальной**, если хотя бы один из составных атрибутов $AB...Z$ совпадает с заголовком r .

Нетривиальные PJD, подразумеваемые возможными ключами, существуют во всех отношениях с арностью, большей двух, первичный ключ которых не совпадает с заголовком отношения.

Примеры:

- СЛУЖ_ПРО_ЗАДАН (СЛУ_N - первичный ключ)
(служащий может работать только в одном проекте и над одним заданием)
то \exists PJD $*(\{СЛУ_N, ПРО_N\}, \{СЛУ_N, СЛУ_ЗАДАН\})$
- $r(ABCD)$ (A – первичный ключ $\Rightarrow A \rightarrow B, A \rightarrow C, A \rightarrow D \Rightarrow PR(AB), PR(AC), PR(AD)$)
то \exists PJD $*(\{AB\}, \{AC\}, \{AD\})$

Нетривиальные PJD, подразумеваемые возможными ключами неинтересны с точки зрения проектирования базы данных, поскольку они не порождают anomalies обновления.

Пятая нормальная форма

Отношение r находится в **пятой нормальной форме**, или в **нормальной форме проекции/соединения (5NF, PJ/NF, Project-Join Normal Form)** в том и только в том случае, когда каждая нетривиальная PJD в r подразумевается возможными ключами r .

5NF является «окончательной» нормальной формой, которой можно достичь в процессе нормализации на основе проекций.

У отношения, находящегося в 5NF, отсутствуют аномалии обновлений, которые можно было бы устранить путем его декомпозиции.

Дальнейшая нормализация отношения в 5NF бессмысленна.

Актуальность процесса нормализации

Первый аспект: Реляционная модель данных != SQL модель данных

Подавляющее большинство современных баз данных и средств управления ими опирается на модель данных SQL, которая не тождественна реляционной модели.

(например, таблица модели SQL может содержать мультимножества строк и не иметь ключа).

Но если потребовать от проектируемой SQL-базы данных наличия хотя бы одного возможного ключа для каждой таблицы целевой базы данных, то все методы нормализации становятся применимы.

Таким образом, SQL-ориентированную базу данных можно проектировать как реляционную базу данных, если должным образом ограничить используемые средства модели данных SQL.

В дальнейшем под «**реляционными**» базами данных будут пониматься именно SQL-ориентированные базы данных, не противоречащие требованиям реляционной модели данных.

Актуальность процесса нормализации

Второй аспект: актуальность «нормализации»

Основные цели процесса нормализации:

- избежать избыточности хранения данных;
- устранить аномалии обновления отношений.

Хорошо нормализованные реляционные базы данных в значительной степени способствуют росту эффективности приложений. И на первых шагах развития это было особенно актуально.

Теория реляционных баз данных и методы их проектирования активно развиваются уже более 25 лет. И ситуация и в области аппаратного и в области программного обеспечения не стоит на месте.

Однако и сейчас наиболее частой схемой использования СУБД остается работа в информационных системах оперативной обработки транзакций (On-Line Transaction Processing – OLTP).

Информационные системы оперативной обработки транзакций On-Line Transaction Processing – OLTP

операционные банковские системы,
системы учета,
системы заказа и резервирования,
и многие другие.

Общая характеристика OLTP систем

- количество транзакций очень велико,
- выполняются транзакции параллельно,
- сложность транзакций низкая,
- необходима высокая эффективность отката транзакции,
- транзакции на вставку/удаления/модификацию превышают запросы на выборку,
- большая часть планируемых запросов известна еще на стадии проектирования.

Вывод:

чем выше уровень нормализации данных в OLTP-приложениях, тем они быстрее и надежнее. Возможны отступления для обеспечения эффективности заранее известных запросов, которые критичны для работы приложений.

Системы категории оперативной аналитической обработки On-Line Analytical Processing – OLAP

OLAP - обобщенный термин, характеризующий принципы построения систем, предназначенных для нахождения зависимостей между данными, для проведения динамического анализа данных и тому подобных задач.

К **OLAP** системам можно отнести системы:

- поддержки принятия решений – Decision Support System (DSS),
- хранилища данных – Data Warehouse,
- системы интеллектуального анализа данных – Data Mining,
- ...

OLAP-приложения оперируют с большими массивами данных и характеризуются следующими признаками:

- добавление в систему новых данных происходит относительно редко крупными блоками;
- данные, добавленные в систему, как правило, никогда не удаляются;
- перед загрузкой данные могут проходить подготовительную обработку;
- запросы к системе являются нерегламентированными и достаточно сложными;
- скорость выполнения запросов важна, но не столь критична как в OLTP системах.

Базы данных OLAP-приложений обычно представлены в виде **гиперкубов** - конструкций, каждое измерение которой представляет собой некоторую характеристику данные, а в ячейках самого гиперкуба хранятся значения этих данных.

Физически гиперкуб может быть построен на основе специальной многомерной модели данных – **Multidimensional OLAP (MOLAP)** или средствами реляционной модели данных – **Relational OLAP (ROLAP)**.

В системах OLAP, использующих реляционную модель данных, данные целесообразно хранить в виде слабо нормализованных отношений, содержащих заранее вычисленные основные итоговые данные. Аномалии обновления мало влияют на работу систем.

Технологически правильным считается для OLAP систем поддерживать отдельную базу данных. Источником данных для такой системы обычно служат различные OLTP системы.

Эффективность работы таких систем обычно сильно зависит от того насколько правильно они спроектированы. Обычно первым шагом проектирования проводят процедуру нормализации схемы системы до максимально возможного уровня нормализации, а вторым шагом проводят процесс денормализации из соображения повышения эффективности выполнения запросов.