

МОДИФИКАЦИЯ АЛГОРИТМОВ НОРМАЛИЗАЦИИ ТАБЛИЦ БАЗ ДАННЫХ ДЛЯ ПРИМЕНЕНИЯ В СИСТЕМАХ BIG DATA



А.А. Карпук

*Профессор кафедры программного обеспечения сетей телекоммуникаций Белорусской государственной академии связи, кандидат технических наук, доцент
a_karpuik@mail.ru*



Л.С. Лазута

*Аспирант кафедры телекоммуникационных систем Белорусской государственной академии связи, магистр
lenya.lazuta@mail.ru*

А.А. Карпук

Окончил Белорусский государственный университет, автор более 260 опубликованных научных трудов, включая 3 монографии. Область научных интересов связана с моделированием и оптимизацией сложных систем, проектированием баз данных и хранилищ данных, оценкой качества радиосвязи и оптимизацией присвоения радиочастот в радиосетях и радиолиниях.

Л.С. Лазута

Окончил Белорусскую государственную академию связи, имеет 8 опубликованных научных трудов. Область научных интересов связана с моделированием и оптимизацией сложных систем, проектированием баз данных и хранилищ данных.

Аннотация. Рассмотрены классические алгоритмы нормализации таблиц в реляционных базах данных. Показаны недостатки алгоритмов, ограничивающие их применение для нормализации таблиц в системах Big Data. Предложены модификации алгоритмов Делобеля-Кейси, Бернштейна и Ислура, которые имеют полиномиальную сложность и могут применяться в системах Big Data для автоматического приведения таблиц данных к третьей нормальной форме.

Ключевые слова: большие данные, функциональная зависимость, третья нормальная форма, алгоритмы нормализации таблиц, минимальное покрытие, замыкание атрибутов, ключ таблицы.

Введение. Современные системы Big Data могут иметь архитектуру реляционного хранилища данных (Relational Data Warehouse, RDW), озера данных (Lake Data), современного хранилища данных (Modern Data Warehouse, MDW), фабрики данных (Data Fabric), озерного хранилища данных (Data Lakehouse), сетки данных (Data Mesh) [1]. Во всех перечисленных архитектурах, кроме архитектуры Lake Data, на одном или нескольких этапах обработки данных решаются задачи очистки данных, предварительной обработки данных и нормализации данных в виде приведения данных к третьей нормальной форме (3НФ). Возможные подходы к решению задач очистки и предварительной обработки данных рассмотрены в работах [2, 3]. Для решения задачи нормализации данных требуется знание функциональных зависимостей (ФЗ), существующих между данными. В работе [4] рассматривались методы выделения и алгоритмы поиска ФЗ между данными в системах Big Data. В настоящей работе рассмотрены классические алгоритмы приведения таблиц реляционной базы данных к 3НФ Делобеля-Кейси, Бернштейна, Ислура и Неклюдовой-Цаленко, показаны недостатки алгоритмов, ограничивающие их применение для нормализации таблиц данных в системах Big Data, и предложены модификации алгоритмов Делобеля-Кейси, Бернштейна и Ислура, которые могут применяться в системах Big Data для автоматического приведения таблиц данных к 3НФ.

Классические алгоритмы нормализации таблиц баз данных. Пусть $X \subseteq A$ – подмножество атрибутов (признаков) системы Big Data, ZA – некоторый атрибут. Говорят, что существует ФЗ XZ , если любой комбинации значений атрибутов из X всегда соответствует единственное значение атрибута Z . Очевидно, что из $Z \in X$ следует, что XZ . Такая ФЗ, в которой зависимый атрибут входит в состав левой части ФЗ, называется тривиальной. В дальнейшем будем рассматривать только нетривиальные ФЗ между атрибутами. Структура ФЗ на множестве атрибутов удовлетворяет аксиомам Армстронга [5]. Для задания структуры ФЗ, отличающейся от тривиальной, требуется постулировать конечное множество ФЗ $F = \{F_j = X_j Y_j | X_j, Y_j A, j=1, m\}$, которое называется системой образующих структуры ФЗ на A . Структуру ФЗ, заданную системой образующих F , будем обозначать $S(F)$. Замыканием множества XA относительно структуры ФЗ $S(F)$ называется множество $X+(F)A$, такое, что для любого YA из XY следует $YX+(F)$.

Структуры ФЗ $S(F_1)$ и $S(F_2)$ на множестве A с системами образующих $F_1 = \{X_i 1 Y_i 1 | X_i 1 \subseteq A, Y_i 1 \subseteq A, i=1, m_1\}$ и $F_2 = \{X_j 2 Y_j 2 | X_j 2 \subseteq A, Y_j 2 \subseteq A, j=1, m_2\}$ соответственно называются эквивалентными, если для любого $X \subseteq A$ имеет место равенство $X+(F_1) = X+(F_2)$. В работе [6] было доказано, что структуры ФЗ $S(F_1)$ и $S(F_2)$ на множестве A эквивалентны тогда и только тогда, когда для всех $i=1, m_1$ выполняются условия $X_i 1+(F_1) = X_i 1+(F_2)$ и для всех $j=1, m_2$ выполняются условия $X_j 2+(F_1) = X_j 2+(F_2)$.

Система образующих $E = \{H_j T_j | H_j, T_j A, j=1, m\}$ структуры ФЗ $S(E)$ на множестве A называется минимальным покрытием (minimal coverage), если выполняются следующие условия:

для любого $j=1, m$ и любого $X T_j$ система образующих E' , полученная из E путем замены ФЗ $H_j T_j$ на ФЗ $H_j T_j \setminus X$, задает структуру ФЗ $S(E')$, не эквивалентную структуре ФЗ $S(E)$;

для любого $j=1, m$ и любого $Y H_j$ система образующих E' , полученная из E путем замены ФЗ $H_j T_j$ на ФЗ $H_j \setminus Y T_j$, задает структуру ФЗ $S(E')$, не эквивалентную структуре ФЗ $S(E)$.

В русскоязычной литературе минимальное покрытие иногда называют элементарным базисом. Классический алгоритм Делобеля-Кейси [7] состоит из следующих шагов.

1. Разложить ФЗ так, чтобы в правой части каждой ФЗ был только один атрибут. Каждая ФЗ $X_j Y_j$ из F , где $Y_j = \{A_{j1}, A_{j2}, \dots\}$, заменяется на несколько ФЗ с одним атрибутом в правой части $X_j A_{j1}, X_j A_{j2}, \dots$.

2. Удалить избыточные атрибуты в левых частях полученных ФЗ. Атрибут VX_j является избыточным в левой части ФЗ $X_j A_{j1}$, если $A_{j1}(X_j \setminus V) +$. При построении замыкания рассматриваемая ФЗ не учитывается.

3. Удалить избыточные ФЗ из полученного множества ФЗ. ФЗ $X_j A_{j1}$ является избыточной, если после ее удаления выполняется условие $A_{j1}(X_j) +$.

4. Все полученные ФЗ с одинаковыми левыми частями собираются в одну ФЗ, в правой части которой находится объединение атрибутов правых частей собранных ФЗ.

5. На основе каждой полученной ФЗ формируется таблица, содержащая все атрибуты, находящиеся в левой и правой части этой ФЗ. Первичным ключом таблицы объявляется множество атрибутов левой части ФЗ.

6. Если среди полученных таблиц найдутся две таблицы, такие что все атрибуты первой таблицы входят в множество атрибутов второй таблицы, то первая таблица удаляется. При этом проверка ФЗ всех неключевых атрибутов первой таблицы от первичного ключа этой таблицы должна выполняться при обработке строк второй таблицы.

7. Построить замыкания первичных ключей полученных таблиц относительно полученной структуры ФЗ. Если ни одно из замыканий не совпадает с множеством атрибутов A , то найти любой ключ этого множества атрибутов и добавить к полученным таблицам в ЗНФ таблицу, состоящую из атрибутов найденного ключа. Для поиска ключа применяется алгоритм последовательного удаления атрибутов из списка всех атрибутов и построения замыкания оставшихся атрибутов. Если удаленный атрибут не входит в полученное замыкание, то он входит в состав ключа и возвращается в список атрибутов.

Делобель и Кейси доказали, что полученные таблицы находятся в ЗНФ, и исходная таблица восстанавливается из полученных таблиц с помощью операций естественного соединения. Легко убедиться, что шаги 1-4 алгоритма Делобеля-Кейси строят минимальное покрытие исходной структуры ФЗ на множестве атрибутов A .

Классический алгоритм Бернштейна [8] состоит из следующих шагов.

1-4. Построить минимальное покрытие структуры ФЗ аналогично шагам 1-4 алгоритма Делобеля-Кейси.

5. Построить замыкание левой части каждой из полученных ФЗ и разбить ФЗ на группы таким образом, чтобы в одной группе были все ФЗ с совпадающими замыканиями левых частей.

6. На основе каждой полученной группы ФЗ формируется таблица, содержащая все атрибуты, находящиеся в левых и правых частях всех ФЗ группы. Левая часть одной из ФЗ группы (любой) объявляется первичным ключом таблицы, левые части остальных ФЗ объявляются уникальными потенциальными ключами таблицы.

7. При формировании таблиц возможен случай, когда в таблице, содержащей более одного ключа, неключевой атрибут, попавший в таблицу из правой части одной ФЗ (ключа), окажется в частичной ФЗ от правой части другой ФЗ, попавшей в таблицу (другого ключа). Такие атрибуты удаляются из таблицы. Для поиска таких атрибутов последовательно рассматриваются неключевые атрибуты таблицы, каждый атрибут временно удаляется из таблицы и из той ФЗ, из которой он попал в таблицу, и строится замыкание оставшихся атрибутов таблицы. Если рассматриваемый атрибут принадлежит этому замыканию, то его действительно надо удалить из таблицы, в противном случае он возвращается в таблицу.

8. Этот шаг аналогичен шагу 6 алгоритма Делобеля-Кейси.

9. Этот шаг аналогичен шагу 7 алгоритма Делобеля-Кейси.

Бернштейн доказал, что полученные таблицы находятся в ЗНФ, и исходная таблица восстанавливается из полученных таблиц в ЗНФ с помощью операций естественного

соединения. Классический алгоритм Ислура [9] совпадает с алгоритмом Бернштейна, за исключением того, что на шаге 9 для поиска ключа всего множества атрибутов используется усовершенствованный алгоритм. Из исходного списка всех атрибутов сразу удаляются все атрибуты, не входящие в левые части исходных ФЗ, а затем оставшиеся атрибуты проверяются на вхождение в ключ таким же образом, как в алгоритмах Делобеля-Кейси и Бернштейна.

На первых шагах классического алгоритма Неклюдовой–Цаленко [10] также строится минимальное покрытие, затем находятся все ключи таблицы, и множество всех ключевых атрибутов таблицы рассматривается как полная подструктура структуры ФЗ. Биекции всех ключей включаются в систему образующих структуры ФЗ, в результате чего в минимальном покрытии могут появиться транзитивные ФЗ, которые удаляются. Производится расширение подструктуры путем добавления неключевых атрибутов, функционально полно зависящих от ключей. Если все атрибуты таблицы войдут в формируемую подструктуру, то она находится в ЗНФ. В противном случае анализируются ФЗ, оставшиеся в минимальном покрытии. Если в них входят не все атрибуты таблицы, то множество атрибутов этих ФЗ рассматривается как полная подструктура структуры ФЗ, для которой все шаги алгоритма повторяются, иначе применяется алгоритм Бернштейна.

Модификация алгоритмов для применения в системах Big Data. Сразу отметим, что поиск всей ключей таблицы, используемый в алгоритме Неклюдовой–Цаленко, является NP-трудной задачей, поэтому применение алгоритма Неклюдовой–Цаленко для автоматической нормализации таблиц систем Big Data, содержащих сотни атрибутов и сотни исходных ФЗ между атрибутами, невозможно.

В алгоритмах Делобеля-Кейси, Бернштейна и Ислура имеется шаг поглощения меньшей таблицы большей таблицей, если все атрибуты меньшей таблицы входят в состав большей таблицы (шаг 6 в алгоритме Делобеля-Кейси). Это фактически ведет к потере ФЗ минимального покрытия, на которой построена меньшая таблица. В базах данных поддержка таких ФЗ производится средствами СУБД с помощью специально написанных триггеров. Большинство проектировщиков баз данных по этой причине не выполняют шаг 6 алгоритма Делобеля-Кейси и шаг 8 алгоритмов Бернштейна и Ислура. Лучше иметь в базе данных лишнюю таблицу, чем решать проблему с поддержкой ФЗ между атрибутами. Мы считаем, что при использовании алгоритмов Делобеля-Кейси, Бернштейна и Ислура для автоматической нормализации таблиц в системах Big Data этот шаг следует исключить из алгоритмов.

Во всех рассмотренных алгоритмах при построении минимального покрытия исходной структуры ФЗ сначала все ФЗ раскладываются на ФЗ, содержащие один атрибут в правой части, а после удаления избыточных атрибутов в левых частях и избыточных ФЗ снова собираются в ФЗ с одинаковыми левыми частями. Применяемый алгоритм удаления избыточных атрибутов в левых частях и избыточных ФЗ основан на том, что в правой части ФЗ находится один атрибут. В работе [11] Карпук предложил алгоритм построения минимального покрытия структуры ФЗ, в котором не требуется раскладывать ФЗ на ФЗ, содержащие один атрибут в правой части, а затем собирать ФЗ с одинаковыми левыми частями. Этот алгоритм состоит из следующих шагов.

1. Удалить избыточные атрибуты из левых частей ФЗ из F. Атрибут VX_j называется избыточным в X_j , если $V(X_j \setminus V) + (F')$, где через F' обозначена система образующих структуры ФЗ, полученная из F путем удаления ФЗ $X_j Y_j$.

2. Если в результате получим в F две или более ФЗ с одинаковыми левыми частями, то объединить их в одну ФЗ.

3. Удалить избыточные атрибуты из правых частей ФЗ из F. Атрибут VY_j называется избыточным в Y_j , если $V(X_j) + (F')$, где через F' обозначена система образующих структуры ФЗ, полученная из F путем замены ФЗ $X_j Y_j$ на $X_j Y \setminus V_j$.

4. Удалить из F ФЗ с пустыми правыми частями.

Карпук доказал, что описанный алгоритм строит минимальное покрытие структуры ФЗ, совпадающее с минимальным покрытием, построенным по шагам 1-4 алгоритма Делобеля-Кейси. Отказ от разложения ФЗ по одному атрибуту в правой части позволяет избежать увеличения мощности множества ФЗ и приводит к снижению вычислительной сложности построения минимального покрытия с $O(m^2k^2ln)$ до $O(m^2(1+k)n)$, где n – количество атрибутов, m – количество ФЗ, k – среднее количество атрибутов в правой части ФЗ, l – среднее количество атрибутов в левой части ФЗ. Предложенный подход обеспечивает асимптотическое ускорение, особенно заметное при наличии ФЗ с большой мощностью правой части, что характерно для систем Big Data. Еще более заметно преимущество алгоритма Карпука для построения минимального покрытия структуры ФЗ при оценке объема оперативной памяти, необходимой для реализации алгоритма. Для реализации алгоритма Делобеля-Кейси объем требуемой оперативной памяти оценивается величиной $O(mkl)$, а для реализации алгоритма Карпука величиной $O_m(k+1)$. Таким образом, при использовании алгоритмов Делобеля-Кейси, Бернштейна и Ислура для автоматической нормализации таблиц в системах Big Data вместо шагов 1-4 для построения минимального покрытия структуры ФЗ следует использовать алгоритм Карпука.

В таблицах систем Big Data количество атрибутов, не входящих в левые части ФЗ, может быть очень большим, поэтому при использовании алгоритмов Делобеля-Кейси, Бернштейна и Ислура для автоматической нормализации таблиц в системах Big Data для поиска ключа всего множества атрибутов следует использовать алгоритм Ислура.

Список литературы

- [1] Serra James. Deciphering Data Architectures. Choosing Between a Modern Data Warehouse, Data Fabric, Data Lakehouse, and Data Mesh. O'Reilly 2024. – 252 p.
- [2] Лазута Л.С. Методы и алгоритмы очистки и предварительной обработки данных в хранилищах данных // Новые информационные технологии в телекоммуникациях и почтовой связи: материалы XXV Междунар. науч.-техн. конф., 13-14 мая 2025 г. / редкол.: А.О. Зеневич [и др.]. – Минск: Белорусская государственная академия связи, 2025. – С. 153–154.
- [3] Lazuta L.S., Karpuk A.A. Cleaning and Preprocessing of Data in Data Warehouses // International Journal of Engineering Research and Development. – 2025. – Vol. 21, Issue 9. – P. 263–269.
- [4] Карпук А.А., Лазута Л.С. Выделение и поиск функциональных зависимостей между атрибутами в системах Big Data // 11-ая Международная научно-техническая конференция «BIG DATA and Advanced Analytics», 23-24 апреля 2025 г. БГУИР, Минск, Беларусь. – С. 397–405.
- [5] Armstrong W.W. Dependency Structure of Data Base Relationships // Proc. IFIP Congress. – Geneva, Switzerland, 1974. – P. 580–583.
- [6] Карпук А.А. Выбор элементарного базиса структуры функциональных зависимостей при проектировании базы данных // Вопросы радиоэлектроники. Сер. ОВР. – 1983. – Вып. 6. – С. 38–41.
- [7] Delobel C., Casey R.G. Decomposition of a data base and the theory Boolean switching functions // IBM J. Res. And Dev. – 1973. – Vol. 17, No 5. – P. 374–386.
- [8] Bernstein P.A. Synthesizing third normal form relations from functional dependencies // ASM Transactions on Database Systems. – 1976. – Vol. 1, No 4. – P. 277–298.
- [9] Isloor S.S. An algorithm with logical simplicity for designing third normal form relational database schema from functional dependencies // Proc. of Int. Conf. on DBMSs (ICMOD 78). – Fast Milan, Italy, 1978. – P. 31–50.
- [10] Неклюдова Е.А., Цаленко М.Ш. Синтез логической схемы реляционной базы данных // Программирование. – 1979. – № 6. – С. 58–68.
- [11] Карпук А.А. О построении элементарного базиса системы функциональных зависимостей в базе данных // Информационные технологии и программные средства: проектирование, разработка и применение: Сб. научн. ст. – Гродно: ГрГУ, 2011. – С. 185–190.

Авторский вклад

Карпук Анатолий Алексеевич – руководство исследованием алгоритмов автоматической нормализации таблиц в таблицах систем Big Data, постановка задач исследования, оценка результатов исследования и определение направлений дальнейшей работы по модификации алгоритмов.

Лазута Леонид Сергеевич – анализ алгоритмов автоматической нормализации таблиц в таблицах систем Big Data, определение области применения алгоритмов, разработка предложений по модификации алгоритмов.

MODIFICATION OF DATABASE TABLE NORMALIZATION ALGORITHMS FOR USE IN BIG DATA SYSTEMS

A.A. Karpuk

*Professor, Department
of Telecommunication Network Software,
Belarusian State Academy of Communications,
PhD of Technical sciences, Associate Professor*

L.S. Lazuta

*Postgraduate student,
Department of Telecommunication Systems,
Belarusian State Academy of Communications,
Master's Degree*

Abstract. Classic algorithms for normalizing tables in relational databases are examined. The algorithms' shortcomings that limit their use for normalizing tables in Big Data systems are highlighted. Modifications to the Delobel-Casey, Bernstein, and Isloor algorithms are proposed. These algorithms have polynomial complexity and can be used in Big Data systems to automatically convert data tables to third normal form.

Keywords: big data, functional dependency, third normal form, table normalization algorithms, minimum coverage, attribute closure, table key.