Санкт-Петербургский государственный университет Программная инженерия

Морозко Иван Дмитриевич

Реализация алгоритма поиска запрещающих ограничений в платформе Desbordante

Отчёт по учебной практике

Научный руководитель: ассистент кафедры ИАС Чернышев Γ . А.

Оглавление

В	ведение	3				
1.	. Постановка задачи					
2.	Обзор					
	2.1. Базовые определения	. 6				
	2.2. Обзор существующих алгоритмов	. 7				
	2.3. Платформа Desbordante	. 9				
3.	Алгоритм FastADC	10				
	3.1. Пространство предикатов	. 10				
	3.2. Набор подсказок и свидетельств	. 11				
	3.3. Приближённые покрытия набора свидетельств	. 13				
	3.4. Алгоритм АЕІ	. 15				
4.	Реализация					
	4.1. Пространство предикатов	. 17				
	4.2. Индекс списка позиций (PLI)	. 17				
5 .	Заключение	19				
Cı	писок литературы	20				

Введение

В последние годы все большее значение приобретают методы обработки и анализа данных. К сожалению, данные, основанные на экспериментах из реального мира, зачастую оказываются неидеальными. В них могут содержаться ошибки, такие как пропущенные значения, неправильные значения и т. п. Недостоверная информация может привести к принятию неверного решения, поэтому очистке данных в современном мире уделяется особое значение [16].

В этом вопросе может помочь профилирование данных — процесс извлечения метаданных из объёма данных. Методов профилирования, представляющий собой поиск некоторой формально описанной закономерности или правила будем называть примитивами. Например, одними из самых известных примитивов являются функциональные зависимости (ФЗ) [11]. Благодаря профилированию данных можно узнать внутренние закономерности в данных, и, благодаря ним, сделать выводы о возможных ошибках [7].

Одним из таких примитивов являются запрещающие ограничения (ЗО). Данный примитив обобщает и расширяет существующий класс ограничений для очистки данных [5]. Каждое ЗО представляет собой отношение между предикатами, указывающее, какие комбинации значений атрибутов являются несогласованными. Автоматических поиск таких закономерностей является трудоёмкой задачей [1], однако, ЗО обладают рядом преимуществ, по сравнению с другими примитивам. Во первых, ЗО позволяют выражать ограничения, не способные быть выраженными другими примитивами. Во вторых, в недавних исследованиях, ЗО использовались в качестве стандартного языка для описания ограничений данных [6,9,12]. Например, HOLOCLEAN [12] — это инструмент для очистки данных, входными данными которого являются наборы приближённых ЗО.

Существует множество алгоритмов поиска ЗО. Более подробно они будут рассмотрены далее, однако среди них выделяется алгоритм Fast-

АDC [10]. Он позволяет находить приближённые запрещающие ограничения на порядок быстрее своих конкурентов. Данный алгоритм полностью написан на языке Java, и представляет собой отдельный проект, не включённый ни в какие известные профилировщики данных, такие как Metanome¹, вследствие чего выявляется ряд проблем. Первая проблема — проблема производительности. Исследование [8] показывает, что производительности Java-приложений не всегда достаточно для вычислений, проводящихся на больших объемах данных. Вторая проблема — удобство конечных пользователей. Зачастую инструментами очистки и профилировки данных пользуются аналитики, далёкие от мира разработки программного обеспечения. В таком случае, необходимость самостоятельной сборки и конфигурации алгоритма может стать препятствием на пути удобного использования.

Для преодоления упомянутых проблем был создан Desbordante², открытая платформа для профилирования данных, написанная на C++. Данный инструмент не поддерживает запрещающие ограничения. Поэтому было решено реализовать и интегрировать в Desbordante алгоритм поиска приближенных запрещающих ограничений — FastADC. Это предоставит открытую, высокопроизводительную и удобную для пользователей реализацию данного алгоритма.

¹hpi.de/naumann/projects/data-profiling-and-analytics/metanome-data-profiling.html

²https://desbordante.unidata-platform.ru/

1 Постановка задачи

Целью данной учебной практики является реализация алгоритма FastADC поиска приближенных запрещающих ограничений в платформе Desbordante. Для её достижения были поставлены следующие задачи:

- 1. Провести обзор предметной области
- 2. Провести сравнение существующих алгоритмов поиска ЗО
- 3. Провести обзор алгоритма FastADC
- 4. Реализовать блоки построения пространства предикатов и индекса списка позиций

2 Обзор

2.1 Базовые определения

Опишем основные понятия, использующиеся для описания запрещающих ограничений и алгоритма их поиска.

Используются стандартные определения из реляционной модели данных: R — отношение; $r \in R$ — атрибут; $s, \ t \in r$ — кортежи.

Определение 1 (Предикат). Запрещающие ограничения определяются на основе предикатов. Каждый предикат p имеет форму $t.A_i$ ор $s.A_j$, где $t, s \in r$, $t \neq s$ (различные кортежи), $A_i, A_j \in R$, и $op \in \{<, \leq, >, \geq, =, \neq\}$. Сравнение проводится по одному и тому же атрибуту, если i = j. В противном случае, оно проводится между двумя атрибутами.

Определение 2 (Запрещающие ограничения). ЗО ϕ на атрибуте r определяется следующим образом: $\forall t, s \in r, \neg (p_1 \land \ldots \land p_m)$, где p_1, \ldots, p_m являются предикатами для кортежей t, s. Мы говорим, что ϕ удерживается на r, если для каждой пары кортежей $(t, s) \in r$ по крайней мере один из p_1, \ldots, p_m не удовлетворяется. Пара кортежей (t, s) нарушает ϕ , если пара удовлетворяет каждому p_i : $i \in [1, m]$.

Определение 3 (Мера ошибки). Значение меры ошибки g_1 для данного 3О ϕ на r, обозначаемое как $g_1(\phi,r)$, определяется как отношение числа пар кортежей, нарушающих ϕ , к общему числу пар кортежей в r. То есть, $g_1(\phi,r) = \frac{|\{(t,s)|(t,s)\in r^2\land (t,s) \text{ нарушает }\phi\}|}{|r|^2-|r|}$.

Определение 4 (Приближённое запрещающее ограничение). При заданном пороге ошибки ϵ , ϕ является приближённым ЗО на r, если и только если $q_1(\phi, r) \leq \epsilon$.

Определение 5 (Минимальное запрещающее ограничение). ЗО ϕ является минимальным, если не существует отличного от ϕ ограничения ϕ' , такого что (a) множество предикатов ϕ' является строгим подмножеством множества предикатов ϕ , и (b) ϕ' является приближённым ЗО на r.

A	В	\mathbf{C}
1	2	фрукт
2	1	овощ
2	4	OBOIII

Таблица 1: Пример таблицы.

Приведём несколько примеров 30, удерживающихся на таблице 1.

Примитив (а) функциональная зависимость $A \to C$ может показать, что любые два кортежа, совпадающие на атрибуте A, совпадают на атрибуте C. Примитив (b) уникальная комбинация колонок (Unique Column Combination, UCC) может показать, что проекция по атрибутам A и B не имеет не имеет одинаковых кортежей. Оба этих примитива можно выразить в терминах 3O:

(a)
$$\forall t, s \in r, \neg (t.A = s.A \land t.C \neq s.C)$$

(b)
$$\forall t, s \in r, \neg (t.A = s.A \land t.B = s.B)$$

(c)
$$\forall t, s \in r, \neg (t.A > s.B \land t.B \le s.B)$$

Так же, ЗО (c) показывает, что возможность использования не только операторов равенства/неравенства, но и числового сравнения, позволяет выражать и более сложные зависимости между данными.

2.2 Обзор существующих алгоритмов

Разработка алгоритмов поиска обычных и приближённых запрещающих ограничений ведётся с 2013 года с представления алгоритма FastDC [4], и постоянно открываются новые методы. Для оценки существующих алгоритмов поиска ЗО и выбора подходящего для реализации алгоритма в платформе Desbordante, было выделено несколько критериев отбора.

Критерии отбора

• Возможность поиска приближённых 30

Несмотря на то, что нахождение приближённых 3О является значительно более трудоёмкой задачей, именно они используются при очистке данных, в отличие от точных 3О [5,17].

• Возможность поиска гетерогенных 30

Возможность использовать предикаты между различными атрибутами позволит увеличить объём информации, которую мы способны получить из данных.

• Нахождение минимальных 30

Нахождение минимальных запрещающих ограничений предпочтительно, поскольку оно позволяет избежать избыточности. Алгоритмы, находящие все 3О, во-первых, используют больше памяти, а во вторых, затрудняют очистку данных. Если 3О ϕ удерживается на r, и существует 3О ξ , следующее из ϕ , то в проверке ξ при очистке данных нет смысла.

• Наличие open-source реализации

Наличие реализации с открытым кодом упрощает процесс реализации алгоритма на языке C++, и позволяет сравнить полученное решение с оригинальным.

• Эффективность по сравнению с предыдущими

Самый эффективный алгоритм является и самым предпочтительным для реализации.

Оценка рассмотренных алгоритмов представлена на таблице 2. Алгоритмы представлены в порядке их выхода, от самого старого на первой строке, до state-of-the-art на последних.

Таблица 2: Оценка рассмотренных алгоритмов.

Алгоритм	Год	Приближённые	Гетерогенные	Минимальные	Открытый	Эффективнее	Константные
		30	30	30	код	предыдущих	30
FastDC [4]	2013	-	+	+	+		+
A-FastDC [4]	2013	+	+	+	+	-	+
Hydra [3]	2017	-	+	+	+	+	-
BFastDC [14]	2018	+	+	+	-	+	+
DCFinder [15]	2019	+	+	+	+	+	-
ADCMiner [2]	2020	+	+	+	+	-	-
DCFinder+ [13]	2022	+	+	+	+*	+	-
FastADC [10]	2022	+	+	+	+	+	-

Возможность поиска 3О, использующих в качестве одно из аргументов некое константное значение, была исследована, однако она не представляет должного интереса чтобы отметить её как необходимую. Данная колонка была помечена тёмно—зелёным цветом как не влияющая на выбор алгоритма.

Хотя авторами работы DCFinder+ была предоставлена открытая реализация³, она, к сожалению, не поддерживает возможность поиска приближённых 3O, описанную в статье. Таким образом, алгоритмами, удовлетворяющим всем критериям, стали DCFinder и FastADC. Однако последний является state-of-the-art алгоритмом. Показано, что FastADC быстрее DCFinder на порядок [10], вследствие чего и был выбран этот алгоритм.

2.3 Платформа Desbordante

Desbordante — это высокопроизводительный инструмент для профилирования данных с открытым исходным кодом, разработанный на C++ [8]. За счет использования низкоуровневых оптимизаций, таких как векторизация вычислений и специализированные аллокаторы, Desbordante обладает высокой производительностью.

Этот инструмент обладает удобным пользовательским интерфейсом, включая Python-биндинги и frontend-клиент.

В настоящее время Desbordante активно развивается, постоянно пополняясь реализациями новых примитивов.

³https://github.com/eduardopena/fdcd

3 Алгоритм FastADC

На текущий момент все алгоритмы поиска приближённых 30 следуют схеме, состоящей из двух этапов [4,10,13,15]. На первом этапе по данному отношению R строится структура, называемая набором свидетельств (evidence set, HC), а на втором этапе происходит поиск приближённых покрытий HC (подробности в разделе 3.3).

FastADC следует такой же схеме, но предлагает новое решение для построения HC, сначала создавая так называемый набор подсказок (clue set), а затем преобразуя его в HC, и предоставляя новый метод для нахождения приближённых ЗО из набора свидетельств — приближённая инверсия свидетельства (Арргохітаte Evidence Inversion, AEI).

Принцип работы алгоритма FastADC можно рассмотреть на рисунке 1.

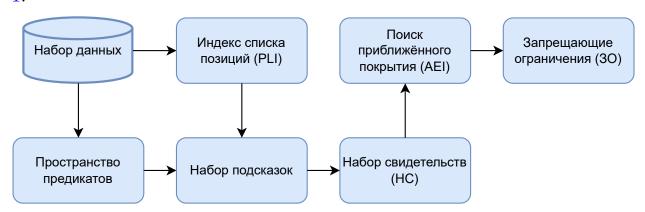


Рис. 1: Схема работы алгоритма FastADC.

3.1 Пространство предикатов

Первым этапом алгоритма FastADC служит построение так называемого пространства предикатов P — множества всех допустимых предикатов на R. Пространство предикатов строится, используя все операторы сравнения для числовых атрибутов, и операторы равенства/неравенства для категориальных. При этом используются только *сравнимые* атрибуты, т.е либо одинаковые, либо атрибуты одинакового типа, между которыми есть не меньше 30% общих значений, что позволяет избегать неинтересных предикатов. Пример пространства предикатов можно увидеть на таблице 3. Данное пространство предикатов построено для таблицы 1, обсуждённой в разделе 2.1.

Таблица 3: Пространство предикатов для таблицы 1.

$p_1: t.A = s.A$	$p_2: t.A \neq s.A$
$p_3: t.A > s.A$	$p_4: t.A < s.A$
$p_5: t.A \ge s.A$	$p_6: t.A \le s.A$
$p_7: t.A = s.B$	$p_8: t.A \neq s.B$
$p_9: t.A > s.B$	$p_{10}: t.A < s.B$
$p_{11}: t.A \ge s.B$	$p_{12}: t.A \le s.B$
$p_{13}: t.B = s.B$	$p_{14}: t.B \neq s.B$
$p_{15}: t.B > s.B$	$p_{16} : t.B < s.B$
$p_{17}: t.B \ge s.B$	$p_{18}: t.B \le s.B$
$p_{19}: t.C = s.C$	$p_{20}: t.C \neq s.C$

3.2 Набор подсказок и свидетельств

Дадим определение набора свидетельств.

Определение 6 (Свидетельство). Свидетельство evi(t,s) для пары кортежей $(t,s) \in r^2$ представляет собой множество предикатов $p \subset P$, удовлетворяющихся на паре кортежей (t,s).

Определение 7 (Набор свидетельств). НС определяется как множество свидетельств по всем парам кортежей: $\{evi(t,s):(t,s)\in r^2\}$.

Обычно НС реализуется с помощью битовых шкал (bit set), что позволяет эффективно выполнять алгоритм поиска приближённых покрытий, описанный в разделе 3.3. Поскольку $evi(t,s) = \bigcup_{A,B \in R} \{\text{предикаты } p: p = (t.A \ op \ s.B)\}$, где $op \in \{=, \neq\}$, либо $\in \{=, \neq, <, >, \geq, \leq\}$, то для кодирования предиката для каждой пары категориальных атрибутов используется два бита, а для пары числовых атрибутов — шесть [14,15].

Однако построение HC из пространства предикатов с данным подходом может быть затруднительно. Как будет показано далее, FastADC строит HC гораздо быстрее существующих алгоритмов, ввиду введения новой структуры — набора подсказок. Подсказка clue(t,s) кодирует отношения между кортежами сравнимых пар атрибутов более сжато, чем набор свидетельств, требует меньше памяти для хранения, и в то же время биективна evi(t,s).

clue(t,s), как и evi(t,s), тоже описывает все предикаты, удовлетворяющиеся на паре кортежей (t,s), однако кодирует **одним** битом предикаты для категориальной пары атрибутов: 0, если $t.A \neq s.B$ (или $t.A \neq s.A$), или 1, если t.A = s.B (или t.A = s.A); и **двумя** битами для численных атрибутов: 00, если t.A < s.B (или t.A < s.A), 01, если t.A = s.B (или t.A = s.A), или 10, если t.A > s.B (или t.A > s.A).

Набор подсказок строится следующим образом. Изначально каждой паре кортежей присваивается начальная подсказка c_0 , после чего подсказки корректируются для пар кортежей, чьи подсказки отличаются от начальной. Эффективность этого процесса достигается за счёт использования индекса списка позиций (Position List Index, PLI)⁴ отношения r. Такой подход значительно эффективнее стандартного метода, сравнивающего все пары кортежей.

Выбор c_0 влияет на эффективность: меньшее количество бит требует коррекции, если c_0 имеет много общих бит с фактическими подсказками. Исходя из предположения что большинство кортежей имеют различные значения атрибутов [15], c_0 выбирается как битовая шкала со всеми значениями равными нулю, т.е., для пары кортежей (t,s) и пары атрибутов (A,B) предполагаем, что $t.A \neq s.B$ для категориальных пар атрибутов и t.A < s.B для числовых пар.

Например, в таблице 1 сравнимыми парами атрибутов являются пары (A,A),(C,C),(A,B) и (B,B). Тогда подсказка clue(2,3) будет выглядеть следующим образом:

$$clue(2,3) = \underbrace{01}_{AA} \underbrace{1}_{CC} \underbrace{00}_{AB} \underbrace{00}_{BB}$$

что соответствует предикатам $\{p_1, p_{19}, p_{10}, p_{16}\}$, соответственно.

В худшем случае сложность алгоритма квадратична по отношению

⁴Подробнее про PLI написано в разделе 4.2.

 $\kappa |r|$, но на практике многие пары кортежей могут быть пропущены благодаря PLI. Также, скорость данного алгоритма в частности достигается за счёт попадания подсказок в кеш процессора. Каждая подсказка может быть скорректирована несколько раз, при этом при коррекции модифицируется максимум один бит для каждой пары атрибутов: бит 0 может быть изменён только на 1 для категориальных пар атрибутов, а 00 может быть изменён только на 01 или 10 для числовых пар.

Генерация доказательств из подсказок происходит через простое преобразование битов [10]. Сложность этого линейна по отношению к размеру набора подсказок, а мощность этого множества $\ll |r|^2$, поскольку многие пары кортежей могут производить одинаковые подсказки, как было показано авторами работы [10]. Тем самым, что время преобразования незначительно по сравнению с временем создания набора подсказок.

3.3 Приближённые покрытия набора свидетельств

Для начала определим понятие обычного и приближенного покрытия HC.

Определение 8 (Покрытие HC). Покрытие HC — это подмножество X множества P, такое что $\forall e \in \text{HC } X \cap e \neq \emptyset$.

Авторы статьи [4], положившей начало алгоритмам поиска приближённых запрещающих ограничений, ввели и доказали следующую теорему:

Теорема 1 (Допустимость точного 3O). ЗО $\phi = \neg(\overline{p_1} \land \cdots \land \overline{p_m})$ допустимо на r, тогда и только тогда, когда множество $\{p_1, \ldots, p_m\}$ является покрытием HC.

С помощью данной теоремы находятся точные запрещающие ограничения. Авторы той же работы [4] также расширяют это определение для приближённых 3O.

Определение 9 (Функция cnt). Разные пары кортежей могут производить одинаковые свидетельства. Функция cnt(e) для свидетельства e обозначает число пар кортежей (t,s), таких что evi(t,s) = e.

Определение 10 (Приближённое покрытие HC). Приближенное покрытие HC это такое подмножество $X \subset P$, что для всех $e \in$ HC, таких что $X \cap e \neq \emptyset$, $\sum cnt(e) \geq (1 - \epsilon) \times (|r|^2 - |r|)$.

Теорема 2 (Валидность приближённого ЗО). При данном пороге ошибки ϵ , приближённое ЗО $\phi = \neg(\overline{p_1} \land \cdots \land \overline{p_m})$ допустимо на r, тогда и только тогда, когда $\{p_1, \ldots, p_m\}$ является приближённым покрытием НС.

Иными словами, множество инвертированных предикатов ЗО должно пересекаться с «достаточным» количеством свидетельств (не обязательно со всеми).

Прежде чем перейти к методу приближенной инверсии доказательств (AEI), введем некоторые обозначения.

- Для ЗО $\varphi = \neg (p_1 \land \ldots \land p_m)$ и свидетельства $e, \varphi \subseteq e$, означает что набор предикатов φ , т.е., $\{p_1, \ldots, p_m\}$, является подмножеством e. Очевидное следствие: $\varphi \subseteq e(t,s) \Rightarrow \varphi$ нарушается парой кортежей (t,s) (но не обязательно наоборот).
- Говорим, что ЗО φ покрывает e(t,s), если $\varphi \not\subseteq e$, и φ не нарушается парой кортежей (t,s).
- Для двух ЗО φ и φ' мы пишем $\varphi \subseteq \varphi'$, если набор предикатов φ является подмножеством набора предикатов φ' . Очевидное следствие: φ' не является минимальным, если $\exists \varphi : \varphi \subsetneq \varphi'$ и φ действителен.
- $\varphi \cup \{p\}$ обозначает 3О, получаемый добавлением нового предиката $p \kappa \varphi$. Данный процесс называется уточнением φ .

3.4 Алгоритм АЕІ

Основная идея алгоритма AEI выглядит следующим образом. Определятся начальный набор $30~\Sigma$, который постепенно модифицируется, и в итоге становится ответом. Этот набор инициализируется запрещающими ограничениями, состоящими из одного предиката:

 $\Sigma \stackrel{init}{=} \{ \phi : \phi = \neg(p) \ \forall p \in P \}$. На каждой итерации алгоритма обрабатывается одно свидетельство e из HC следующим образом: для каждого $\phi \in \Sigma$ проверяется, покрывает ли оно e. Если ϕ не покрывает свидетельство, алгоритм пытается покрыть ϕ e: к ϕ добавляются новые предикаты из P, которые не входят в e, создавая новые кандидатные 3O. После каждого добавления нового предиката к ϕ проверяется, остается ли полученное 3O минимальным. Если нет, 3O исключается дальнейшего рассмотрения. Итерации повторяются пока каждое свидетельство не будет рассмотрено.

Однако данное описание сможет найти только точные 3О, поскольку при генерации кандидатов приближенных 3О необязательно разрешать все нарушения 3О относительно каждого свидетельства e. Вместо этого должно учитываться количество (cnt(e)) для каждого e при генерации или отсечении кандидатов.

Модификации для нахождения приближённых 30

Каждое кандидатное ЗО ϕ хранится в паре с множеством предикатов cand, которые могут быть использованы для ymoчhehus ϕ . Множество таких пар называется Ψ .

Как и раньше, на каждой итерации рассматривается свидетельство $e \in HC$. Все кандидатные 3О из Ψ делятся на две части в зависимости от того, покрывают ли они e. 3О которые не покрывают e, помещаются в Ψ^- , в то время как остальные остаются в Ψ .

Затем AEI пытается покрыть 3О $\phi \in \Psi^-$, как и ранее описанный алгоритм, но в этот раз сначала с использованием cnt проверяется, может ли этот путь привести к корректному 3О. Вычисляется количество оставшихся свидетельств, которые могут быть покрыты остав-

шимися предикатами, т.е. верхняя граница накопленного количества, и данное число сравнивается с N — суммарным числом свидетельств, которые обязаны быть покрытыми (инициализируется выражением $(1-\epsilon)\times(|r|^2-|r|)$ из определения 10). Если полученное число меньше N, то из cand удаляется e. Если после этого cand представляет собой пустое множество, то ϕ больше невозможно ymovnsmb. Происходит проверка того что ϕ является минимальным и приближённым 3О, и ϕ добавляется в Σ . Происходит рекурсивный вызов, которому передаётся Ψ^- и следующее свидетельство.

В ветви, где e покрыто, собирается отдельное множество предикатов $\{\bigcup_{\{_,cand\}\in\Psi^-} cand \setminus e\}$. Для каждого предиката p из этого множества формируется новое кандидатное ЗО путём ymovinehus $\phi \in \Psi$: $\phi' = \phi \cup \{p\}$, и определяется набор cand' с предикатами из cand, определённых на $\underline{\text{тех же}}$ парах атрибутов что и p. Если такие предикаты существуют, то Ψ обновляется — в него добавляется новая пара $\{\phi', cand'\}$, если, конечно, ϕ' минимальна в Ψ , т.е $\nexists \phi \in \Psi : \phi \subset \phi'$.

Если же cand' представляет собой пустое множество, происходит то же самое что и в первой ветви: проверка того, что ϕ' является минимальным и приближённым 3O; добавление ϕ' в Σ .

Обе ветви пройдены. N уменьшается на cnt(e) и метод рекурсивно вызывается, в этот раз для $\Psi.$

Более подробно алгоритм AEI можно рассмотреть в работе [10].

4 Реализация

Были реализованы необходимые классы по описанию пространства предикатов и построению индекса списка позиций.

Тем самым, реализация содержит два блока алгоритма FastADC, нужных для построения *набора подсказок*.

4.1 Пространство предикатов

Пространство предикатов представляет собой список классов $\Pi pe-du\kappa am$. Для построения класса $\Pi pedu\kappa am$, описывающий предикат $p=t.A\ op\ s.B$, были построены классы $OnepandKonon\kappa u$ и Onepamop.

Алгоритм построения пространства предикатов проходится по каждой паре атрибутов и вычисляет, являются ли атрибуты сравнимыми. Напомним, что *сравнимые* атрибуты, это либо одинаковые атрибуты, либо атрибуты одинакового типа, между которыми есть не меньше 30% общих значений. Были реализованы необходимые функции, возвращающие метрику сравнимости для атрибутов категориальных и числовых типов.

4.2 Индекс списка позиций (PLI)

Каждый PLI строится для атрибута $A \in R$. При этом предварительно значения атрибута хешируются, что позволяет говорить о значении кортежа на атрибуте A как об обычном числе — ключе, не заботясь о настоящем типе атрибута.

Определение 11 (Кластер). Кластером в контексте PLI мы называем пару $\langle k, l \rangle$, где ключ k является значением в атрибуте A, а l — множество всех кортежей, имеющих то же значение k в атрибуте A.

Определение 12 (Индекс списка позиций (PLI)). Для атрибута A, класс PLI является тройкой:

- ullet Множество ключей уникальных значений атрибута A
- Список кластеров
- ullet Отображение ключ o кластер

Для числового атрибута A, кластеры в PLI дополнительно сортируются по ключу k в порядке убывания, что позволит оптимизировать построение набора подсказок путём использования двоичного поиска.

Данный класс был реализован, наряду с методом хеширования значений атрибутов таблицы.

5 Заключение

В результате проделанной работы были выполнены следующие задачи:

- Проведён обзор предметной области
- Выбран набор критериев, по которым проведён обзор алгоритмов поиска 30
- Проведён обзор алгоритма FastADC
- Реализованы блоки пространства предикатов и PLI алгоритма FastADC

Код находится в репозитории 5 .

⁵https://github.com/ol-imorozko/Desbordante/commits/FastADC/

Список литературы

- [1] Abedjan Ziawasch, Golab Lukasz, and Naumann Felix. Profiling relational data: a survey // The VLDB Journal. 2015. Vol. 24, no. 4. P. 557–581. Access mode: https://doi.org/10.1007/s00778-015-0389-y.
- [2] Livshits Ester, Heidari Alireza, Ilyas Ihab F., and Kimelfeld Benny. Approximate Denial Constraints. 2020. 2005.08540.
- [3] Bleifuß Tobias, Kruse Sebastian, and Naumann Felix. Efficient denial constraint discovery with hydra // Proc. VLDB Endow. 2017. Vol. 11, no. 3. P. 311–323. Access mode: https://doi.org/10.14778/3157794.3157800.
- [4] Chu Xu, Ilyas Ihab F., and Papotti Paolo. Discovering denial constraints // Proc. VLDB Endow. 2013. Vol. 6, no. 13. P. 1498—1509. Access mode: https://doi.org/10.14778/2536258.2536262.
- [5] Chu Xu, Ilyas Ihab F., and Papotti Paolo. Holistic data cleaning: Putting violations into context // 2013 IEEE 29th International Conference on Data Engineering (ICDE). 2013. P. 458–469.
- [6] Fagin Ronald. Kimelfeld Benny, Reiss Frederick. and Vansummeren Stijn. Cleaning inconsistencies in information extraction via prioritized repairs // Proceedings of the 33rd ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems. — New York, NY, USA: Association for Computing Machinery. — 2014. — PODS '14. — P. 164–175. — Access mode: https://doi.org/10.1145/2594538.2594540.
- [7] Bohannon Philip, Fan Wenfei, Geerts Floris, Jia Xibei, and Kementsietsidis Anastasios. Conditional Functional Dependencies for Data Cleaning // 2007 IEEE 23rd International Conference on Data Engineering. 2007. P. 746–755.

- [8] Strutovskiy Maxim, Bobrov Nikita, Smirnov Kirill, and Chernishev George. Desbordante: a Framework for Exploring Limits of Dependency Discovery Algorithms // 2021 29th Conference of Open Innovations Association (FRUCT). 2021. P. 344–354.
- [9] Fan Wenfei, Geerts Floris, and Wijsen Jef. Determining the Currency of Data // ACM Trans. Database Syst. 2012. Vol. 37, no. 4. Access mode: https://doi.org/10.1145/2389241.2389244.
- [10] Xiao Renjie, Tan Zijing, Wang Haojin, and Ma Shuai. Fast approximate denial constraint discovery // Proc. VLDB Endow. 2022. Vol. 16, no. 2. P. 269–281. Access mode: https://doi.org/10.14778/3565816.3565828.
- [11] Papenbrock Thorsten, Ehrlich Jens, Marten Jannik, Neubert Tommy, Jan-Peer, Schönberg Martin, Zwiener Jakob, Naumann Felix. Functional dependency discovery: an experimental evaluation of seven algorithms // Proc. VLDB Endow. — 8, 10. — P. 1082–1093. — Access 2015. - Vol.no. https://doi.org/10.14778/2794367.2794377.
- [12] Rekatsinas Theodoros, Chu Xu, Ilyas Ihab F., and Ré Christopher. HoloClean: holistic data repairs with probabilistic inference // Proc. VLDB Endow. 2017. Vol. 10, no. 11. P. 1190–1201. Access mode: https://doi.org/10.14778/3137628.3137631.
- [13] Pena Eduardo H. M., Porto Fabio, and Naumann Felix. Fast Algorithms for Denial Constraint Discovery // Proc. VLDB Endow.— 2022. Vol. 16, no. 4. P. 684–696. Access mode: https://doi.org/10.14778/3574245.3574254.
- [14] Pena Eduardo H. M. and de Almeida Eduardo Cunha. BFASTDC: A Bitwise Algorithm for Mining Denial Constraints // Database and Expert Systems Applications / ed. by Hartmann Sven, Ma Hui, Hameurlain Abdelkader, Pernul Günther, and Wagner Roland R. Cham: Springer International Publishing. 2018. P. 53–68.

- [15] Pena Eduardo H. M., de Almeida Eduardo C., and Naumann Felix. Discovery of approximate (and exact) denial constraints // Proc. VLDB Endow. 2019. Vol. 13, no. 3. P. 266–278. Access mode: https://doi.org/10.14778/3368289.3368293.
- [16] Rahm Erhard and Do Hong Hai. Data Cleaning: Problems and Current Approaches. // IEEE Data(base) Engineering Bulletin. 2000. Vol. 23. P. 3–13. Access mode: https://api.semanticscholar.org/CorpusID:260972099.
- [17] Rezig El Kindi, Ouzzani Mourad, Elmagarmid Ahmed K., Aref Walid G., and Stonebraker Michael. Towards an End-to-End Human-Centric Data Cleaning Framework // Proceedings of the Workshop on Human-In-the-Loop Data Analytics. New York, NY, USA: Association for Computing Machinery. 2019. HILDA '19. Access mode: https://doi.org/10.1145/3328519.3329133.