

Высокопроизводительный поиск функциональных зависимостей в данных

Отчёт по учебной практике

Салью Артур Кристофович
группа 19.Б11-мм

Научный руководитель:
асс. кафедры ИАС Чернышев Г.А.

Введение

Представим привычную нам таблицу с данными.

Пусть два её столбца имеют имена A и B соответственно.

Если каждому элементу столбца A *однозначно* сопоставлен элемент из B , тогда говорят, что соблюдается **функциональная зависимость**

$$FD: A \rightarrow B$$

Пример функциональных зависимостей

Отношение R			
Имя	Фамилия	Марка автомобиля	Цвет автомобиля
Александр	Петров	Desbordauto	Navy blue
Андрей	Серый	Desbordauto	Snakeskin Green
Василиса	Шавел	Fdepus	Ghost White
Дмитрий	Петров	Desbordauto	Navy blue
Василиса	Шавел	FdepuZ	Ghost White

Алгоритмы поиска

- Поиском ФЗ занимаются достаточно давно;
- ФЗ применяются при нормализации БД, поиске ошибок, дубликатов, анализе данных;
- Алгоритмы поиска ресурсоёмки.

Существующее решение

Metanome — первая платформа с алгоритмами поиска.

Metanome написан на Java:

- необходимость тонкой настройки на уровне JVM;
- затруднён анализ производительности;
- высокие затраты памяти (по сравнению с C++).

Решение — реализуем алгоритм на C++ для проекта Desbordante

Цели и задачи

Цель работы — реализовать стартовую версию алгоритма поиска функциональных зависимостей (FDep) на языке программирования C++ для платформы Desbordante и измерить её производительность.

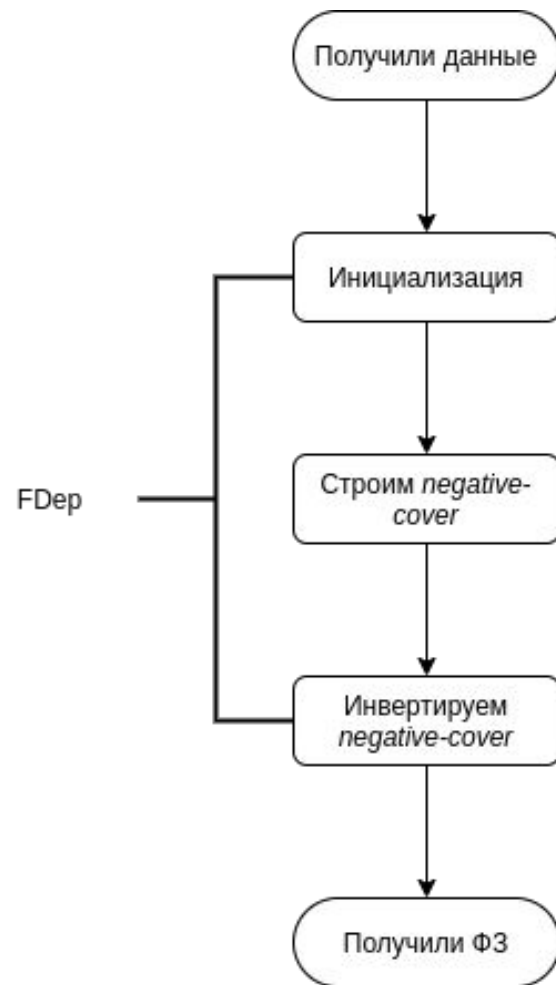
Задачи:

- обзор предметной области;
- идейный обзор алгоритма FDep;
- реализация алгоритма FDep на C++;
- сравнительный анализ производительности алгоритма.

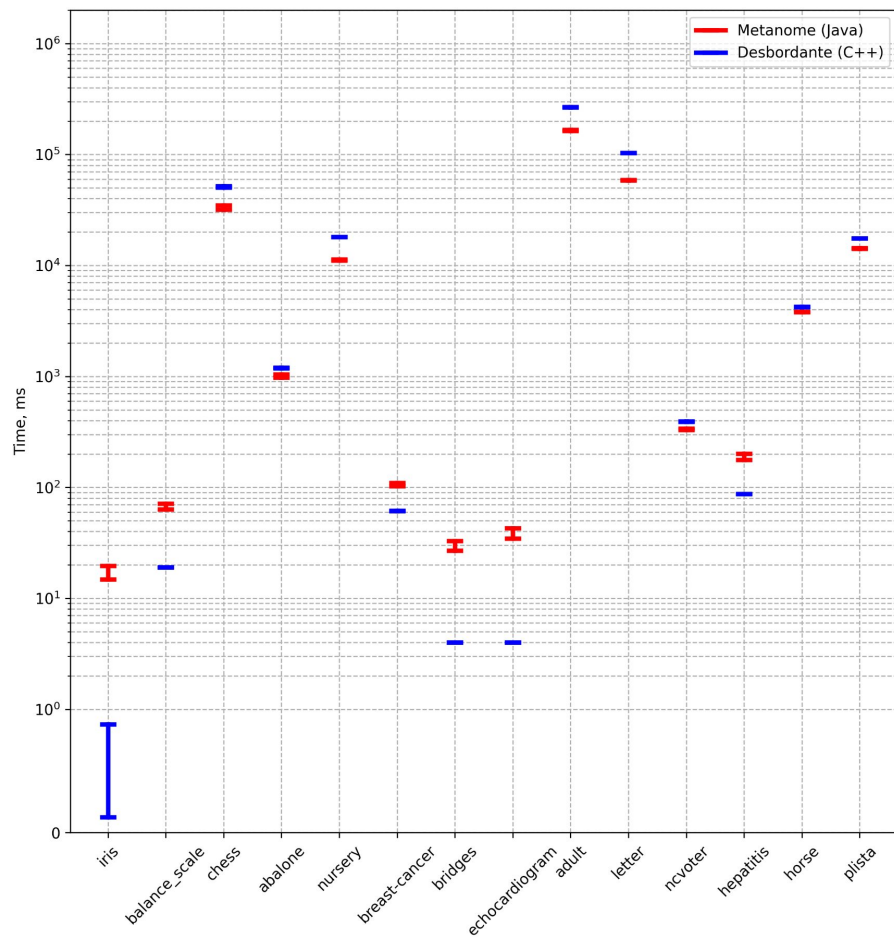
Алгоритм FDep

- Индуктивный алгоритм “от общего к частному”;
- *bottom-up* подход показывает высокую эффективность;
- основная структура — дерево зависимостей.

Блок-схема работы FDep



Сравнение с Metanome



3 датасета - выигрыш (в 1.5 раза)

4 датасета - приблизительно одинаково

5 датасетов - проигрыш (в 1.5 раза)

Заключение

Проведен обзор предметной области:

- обзор области применения функциональных зависимостей;
- обзор алгоритма FDep.

А так же

- реализован алгоритм FDep на языке C++;
- проведен анализ производительности алгоритма.

Дополнительно (пояснение работы FDep)

	A	B	C	D
0	0	0	0	0
1	1	1	0	0
0	0	2	0	2
1	1	2	3	4

Получение NCOVER:

t1, t2: $CD \rightarrow A$, $CD \rightarrow B$

t1, t3: $AC \rightarrow B$, $AC \rightarrow D$

Если продолжить, получим множество ВСЕХ нарушенных ФЗ. Оставим с минимальной LHS.

NCOVER: $\{B \rightarrow A, CD \rightarrow A, AC \rightarrow B, CD \rightarrow B, A \rightarrow C, B \rightarrow C, AC \rightarrow D, B \rightarrow D\}$

Преобразование к positive-cover

На предыдущем шаге получено $NCOVER = \{B \rightarrow A, CD \rightarrow A, AC \rightarrow B, CD \rightarrow B, A \rightarrow C, B \rightarrow C, AC \rightarrow D, B \rightarrow D\}$
 $DEPS = \{\emptyset \rightarrow A, \emptyset \rightarrow B, \emptyset \rightarrow C, \emptyset \rightarrow D\}$ - начальное множество наиболее общих зависимостей.

Во внешнем цикле рассматриваем $ND: B \rightarrow A$.

Из каких зависимостей в $DEPS$ следует ND ? Ответ - $\{\emptyset \rightarrow A\}$

Для каждой зависимости D в этом множестве делаем следующее:

Строим $spec(D, ND) = \{C \rightarrow A, D \rightarrow A\}$ - специализации, из которых не следует ND .

Далее $DEPS = (DEPS / \mathbf{D}) \cup spec(\mathbf{D}, ND) = \{C \rightarrow A, D \rightarrow A, \emptyset \rightarrow B, \emptyset \rightarrow C, \emptyset \rightarrow D\}$

Теперь рассматриваем $CD \rightarrow A$. Эта ФЗ следует из $\{C \rightarrow A, D \rightarrow A\}$

$spec(C \rightarrow A, CD \rightarrow A) = \{BC \rightarrow A\}$

$spec(D \rightarrow A, CD \rightarrow A) = \{BD \rightarrow A\}$

$DEPS = \{BC \rightarrow A, BD \rightarrow A, \emptyset \rightarrow B, \emptyset \rightarrow C, \emptyset \rightarrow D\}$

И так далее... Получим ответ:

POS_COV = $\{BC \rightarrow A, BD \rightarrow A, AD \rightarrow B, AB \rightarrow C, D \rightarrow C, AB \rightarrow D, BC \rightarrow D\}$