

Санкт-Петербургский государственный университет

Математическое обеспечение и администрирование информационных
систем

Системное программирование

Группа 23.М04-мм

Черников Антон Александрович

Расширение возможностей
профилировщика данных Desbordante по
работе с графовыми зависимостями

Отчёт по учебной практике

Научный руководитель:
ассистент кафедры ИАС Г. А. Чернышев

Санкт-Петербург
2024

Оглавление

Введение	3
1. Постановка задачи	5
2. Предварительные сведения	6
3. Обзор	10
4. Алгоритм поиска зависимостей	11
4.1. Входные параметры	11
4.2. Предобработка	13
4.3. Основная работа	15
5. Реализация	19
5.1. Алгоритм	19
5.2. Тестирование	21
5.3. Python bindings	21
5.4. Примеры	22
5.5. Python CLI	24
Заключение	25
Список литературы	26

Введение

Профилирование данных представляет собой распространённую задачу для специалистов, работающих с большими объёмами информации. Этот процесс включает в себя извлечение дополнительной информации о данных, такой как информация об авторе, дата создания или изменения, а также размер занимаемой памяти. Однако помимо этого данные могут содержать неочевидные зависимости и закономерности между элементами, сокрытые в них. В данной работе рассматривается вопрос выявления такого рода информации из данных.

Desbordante¹ — это высокопроизводительный инструмент для профилирования данных с открытым исходным кодом, разработанный группой студентов под руководством Г. А. Чернышева. Проект содержит множество алгоритмов, которые способны обнаруживать различные закономерности в данных, а также предоставляет соответствующие пользовательские интерфейсы для них. В качестве основного языка используется C++, что в целом повышает производительность.

Графовые функциональные зависимости представляют собой естественное обобщение традиционных функциональных зависимостей на структуры данных, такие как графы [4]. Они помогают выявлять несоответствия в базах знаний, находить ошибки, определять спам и управлять блогами в социальных сетях.

Desbordante уже содержит алгоритмы, работающие с графовыми зависимостями, такие как алгоритмы проверки выполнимости существующих графовых зависимостей. Эти алгоритмы получают на вход граф и множество графовых функциональных зависимостей, после чего возвращают только те из них, которые выполнены на данном графе.

Проверка наличия существующей зависимости в графе — очень нужная задача, однако не всегда может быть известно, какой именно структурой обладает интересующая зависимость. Полезным алгоритмом, расширяющим взаимодействие с графовыми зависимостями, является автоматическая генерация графовых зависимостей на основе

¹<https://github.com/Desbordante> (дата обращения 1.05.2024)

входного графа. Авторы понятия графовых зависимостей разработали такой алгоритм и описали его в статье [1], однако ими не был предоставлен код алгоритма. Данная работа направлена на создание первой публичной реализации алгоритма поиска графовых зависимостей. В дальнейшем планируется улучшение и оптимизация реализованного алгоритма.

1. Постановка задачи

Целью данной работы является расширение инструментария Desbordante для работы с графовыми зависимостями.

Для достижения этой цели были поставлены следующие задачи:

- Реализовать алгоритм поиска графовых функциональных зависимостей и произвести тестирование.
- Реализовать возможность запускать алгоритм поиска графовых зависимостей из скриптов, написанных на языке программирования Python.
- Создать скрипты-примеры работы алгоритма поиска графовых зависимостей на языке программирования Python.
- Обеспечить возможность запускать алгоритм поиска графовых зависимостей через консоль путём реализации соответствующей подсистемы.

2. Предварительные сведения

Определение 1 (Функциональная зависимость) *Отношение R удовлетворяет функциональной зависимости $X \rightarrow Y$ (где $X, Y \subset R$) тогда и только тогда, когда для любых кортежей $t_1, t_2 \in R$ выполняется: если $t_1[X] = t_2[X]$, то $t_1[Y] = t_2[Y]$.*

Таблица 1: Данные о студентах и их оценках

ID	Name	Course	Grade
1	Alice	Math	A
2	Bob	Math	B
3	Charlie	Science	A
1	Alice	Science	B
2	Bob	Science	A

Пусть, отношение представлено в виде Таблицы 1. Заметим, что функциональная зависимость $ID \rightarrow Name$ выполнена, так как в этом случае каждый идентификатор студента (ID) уникально определяет его имя (Name). Например, для $ID = 1$ всегда будет $Name = \text{“Alice”}$. В это же время зависимость $Name \rightarrow Course$ не выполнена. Здесь мы видим, что одно и то же имя может соответствовать нескольким курсам. Например, “Alice” изучает как “Math”, так и “Science”. Это означает, что имя не может однозначно определить курс, что делает эту зависимость невыполненной.

Функциональные зависимости могут быть обобщены на графы. Одно из таких обобщений предлагают авторы статьи [4], на котором и основана данная работа. В этой статье определяются и исследуются графовые зависимости, формулируется задача проверки (validation) выполнения зависимостей на графе, а также задачи выполнимости (satisfiability) и импликации (implication) набора зависимостей.

Задачи выполнимости и импликации были более подробно изучены в статье [3], в которой предложены эффективные алгоритмы работы под каждую из них.

Прежде чем рассматривать графовые зависимости, нужно формально определить данные, на которых они определены — графы.

Определение 2 (Граф) *Граф — это структура данных, состоящая из четвёрки (V, E, L, A) , где V — множество вершин; $E \subseteq V \times V$ — множество рёбер; $L : V \cup E \rightarrow \Sigma$ — сюръекция, где Σ — множество меток (алфавит), A — функция, которая сопоставляет каждой вершине список её атрибутов.*

Список атрибутов содержит названия атрибутов и соответствующие этим атрибутам значения. Пусть, $A(u) = (f_1 = c_1, f_2 = c_2, \dots, f_m = c_m)$, $u \in V$, здесь вершина u имеет атрибуты f_i $i = 1, 2, \dots, m$, а число m зависит от конкретной вершины, то есть, у каждой вершины может быть свой набор атрибутов (обычно набор атрибутов зависит от метки вершины). c_i — значение, которое принимает атрибут f_i , обозначение: $u.f_i = c_i$.

В данной работе графы рассматриваются как неориентированные, то есть, $(u, v), (v, u) \in E$ представляют собой один и тот же объект.

Определение 3 (Графовая функциональная зависимость)

GFD (Graph Functional Dependency) — это конструкция $P[X \rightarrow Y]$, где P — паттерн, а X и Y — множества литералов.

В этом определении под паттерном понимается граф, вершины которого однозначно проиндексированы от 0 до $|V| - 1$ для получения доступа к ним, а под литералом — выражение, имеющее вид $i.f = c$ (константный литерал), где i — индекс вершины паттерна, f — атрибут соответствующей вершины, c — константа (значение), или $i.f_i = j.f_j$ (переменный литерал), где i, j — индексы вершин паттерна, f_i, f_j — атрибуты соответствующих вершин.

На Рис. 1 представлен пример графа. Вершины графа имеют метки C или U . В зависимости от метки вершина имеет свой собственный набор атрибутов. В данном примере все вершины имеют одноэлементный список атрибутов. У вершин с меткой C он состоит из элемента *topic*, а у вершин с меткой U — *age_group*. Конкретные значения этих атрибутов указаны рядом с вершинами.

На Рис. 2 продемонстрированы графовые зависимости. Чтобы проверить, выполняется ли GFD, необходимо найти все подграфы графа,

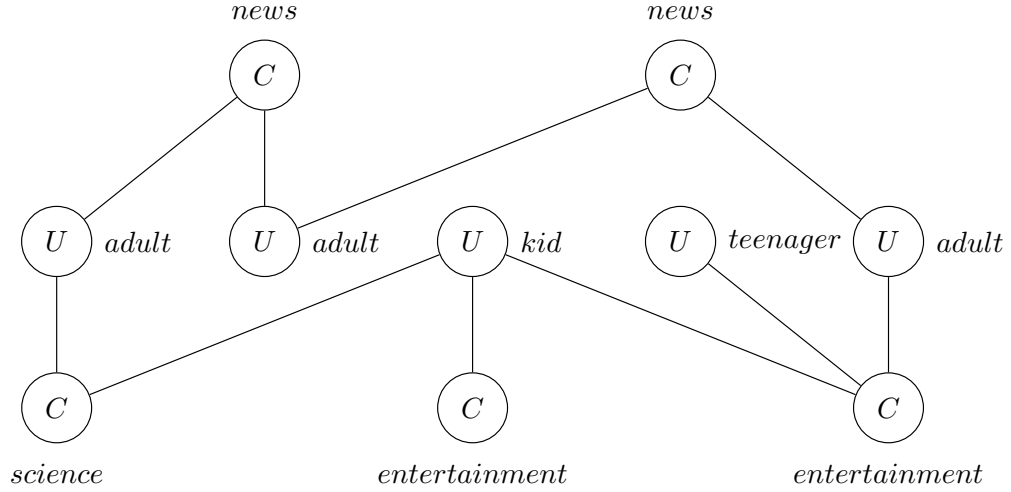


Рис. 1: Связь каналов C (Channel) с пользователями U (User). Атрибуты у вершин с меткой C — $\{topic\}$, с меткой U — $\{age_group\}$.



$$\{0.topic = news\} \rightarrow \{1.age_group = adult\} \quad \{0.topic = news\} \rightarrow \{1.age_group = adult\}$$

(a) Графовая зависимость.

(b) Минимальная графовая зависимость.

Рис. 2: Пример графовых функциональных зависимостей.

изоморфные паттерну, и на каждом вложении проверить выполнимость зависимости литералов, то есть, выполнено ли: если все литералы в левой части выполняются, то все литералы в правой части так же выполняются. Если есть хотя бы одно вложение, на котором зависимость не выполняется, то графовая зависимость не выполняется на всём графе. Если не нашлось ни одного вложения паттерна, то такая зависимость считается тривиально выполненной.

Нетрудно заметить, что в данном примере обе зависимости выполняются на приведённом графе.

Допустим, дан граф G и GFD φ . Примем обозначение: $G \models \varphi$ означает, что графовая зависимость φ выполнена на графе G .

Определение 4 (Минимальная GFD) Пусть $GFD \varphi = P[X \rightarrow Y]$ выполнена на графе G .

φ — минимальная, если $\forall \psi = Q[X \rightarrow Y] : G \models \psi \ P \subset Q$.

Рассмотрим графовые зависимости на Рис. 2. Графовая зависимость, изображённая на Рис. 2b, говорит о том, что в графе все подписчики новостных каналов являются взрослыми. Однако, если мы добавим новую вершину к паттерну, как показано на Рис. 2a, полученная графовая зависимость будет так же выполняться на графе, хотя она несёт в себе избыточную информацию и перегружает понимание сути. Именно по этой причине в данной работе будут находиться только минимальные графовые зависимости.

3. Обзор

Графовые зависимости были предложены авторами статьи [4], в которой также описывают и оценивают алгоритм проверки выполнения набора графовых зависимостей на больших реальных графах.

В результате одной из предыдущих работ [7] был реализован и интегрирован алгоритм проверки графовых функциональных зависимостей в проект Desbordante. Для него были реализованы три версии. Одна из них наивная, необходимая для сравнения с остальными алгоритмами. Вторая является реализацией алгоритма из рассмотренной статьи, а третья — сконструированная улучшенная его версия [5], использующая эффективный алгоритм поиска подграфа CPI [2]. Кроме этого, проект Desbordante был расширен возможностью запускать интегрированный алгоритм проверки графовых зависимостей на Python и через консоль, а также снабжён скриптами-примерами работы этого алгоритма.

Другая предыдущая работа [8] была направлена на обзор новой статьи [1], которая предлагает принципиально другой алгоритм — алгоритм поиска графовых функциональных зависимостей. Их различие с алгоритмом проверки графовых зависимостей в том, что в первом случае перед работой алгоритма пользователю известна вся информация об интересующей зависимости, и задача алгоритма — дать ответ на вопрос выполняется ли данная зависимость на графе. А алгоритм поиска зависимостей позволяет генерировать такие зависимости автоматически. Исходные данные ограничиваются лишь графом, на котором необходимо произвести поиск выполненных зависимостей.

Данная работа направлена на реализацию этого алгоритма поиска графовых зависимостей.

4. Алгоритм поиска зависимостей

Алгоритм поиска графовых функциональных зависимостей, который будет рассматриваться в данной работе, описан в статье [1].

4.1. Входные параметры

Дан граф $G = (V, E, L, A)$. Алгоритм состоит из двух этапов: генерация паттернов и генерация правил литералов. После этого проверяется выполнимость сгенерированных зависимостей.

Сложность заключается в том, что количество подграфов-кандидатов с n вершинами растёт экспоненциально с ростом n . Из-за этого долго проверять существование вложений наивно сгенерированных подграфов. Чтобы решить эту проблему, авторы предлагают использовать несколько эвристик, которые существенно сокращают время работы алгоритма.

Во-первых, помимо графа на вход алгоритму подаётся целое число k , указывающее на то, какое максимальное количество вершин ожидается от паттернов найденных функциональных зависимостей. Если приравнять этот параметр к количеству вершин в графе, то алгоритму придётся перебирать всевозможные подграфы. Однако, зависимости, содержащие паттерны с огромным количеством вершин, в большинстве случаев очень неинформативны, поэтому у пользователя нет необходимости искать зависимости с такими паттернами. Вместо этого он может указать, сколько конкретно вершин ему хотелось бы видеть в итоговом результате.

Во-вторых, алгоритм так же получает пороговое значение σ , которое выражает минимальное значение метрики, условно обозначающей количество вложений паттерна в граф, на которых графовая зависимость выполнена. Можно считать, что это минимальная частота встречаемости паттерна в графе. Если зависимость присутствует в графе в единственном экземпляре или по крайней мере всего в нескольких, то скорее всего пользователю не нужна эта информация. С другой стороны алгоритм является итеративным и использует предыдущие результаты

для следующей итерации. Сокращение результатов на текущей итерации упростит работу следующих шагов. Пользователь может выставить желаемую частоту, чтобы увидеть только популярные зависимости.

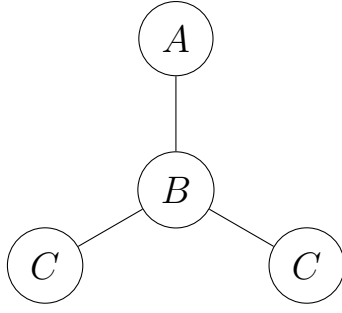
Прежде чем перейти дальше, стоит обсудить вышеупомянутую метрику. Пусть, она будет выражаться функцией $\text{supp}(\varphi, G)$, где φ — графовая зависимость. От этой метрики требуется выполнение двух свойств:

1. Асимптотически она должна выражать частоту встречаемости данной графовой зависимости в графе.
2. Чем больше элементов в паттерне графовой зависимости, тем меньше должно быть значение этой метрики.

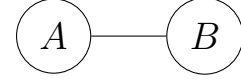
Выполнение первого свойства необходимо по очевидным причинам: от пользователя ожидается, что он будет вводить значение этой метрики исходя из ожиданий популярности интересующих зависимостей. Что касается второго свойства, то оно необходимо для того, чтобы сократить количество генерируемых паттернов на текущей итерации, а, следовательно, и время работы алгоритма.

Если взять за функцию $\text{supp}(\varphi, G)$ количество вложений паттерна графовой зависимости φ в граф G , на которых она выполняется, то второе свойство не будет выполняться. Рассмотрим пример на Рис. 3. Легко увидеть, что паттерн P_1 имеет лишь одно вложение в данный граф G . Добавим теперь к нему новую вершину с меткой C и ребро между вершинами с метками B и C . Полученный паттерн P_2 имеет уже два вложения в тот же граф.

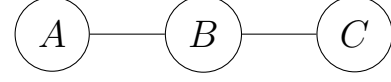
Исходя из вышесказанного, видно, что самое естественное определение функции supp не подходит для алгоритма. Усовершенствуем эту метрику так, чтобы второе свойство выполнялось. Предлагается посчитать для каждой вершины паттерна количество вершин-образов графа, на вложении в которые при этом выполняется данная зависимость. После этого взять минимальное значение среди посчитанных за значение метрики. Нетрудно заметить, что теперь второе свойство будет



(a) Граф G .



(b) Паттерн P_1 ; вложений в G : 1.



(c) Паттерн P_2 ; вложений в G : 2.

Рис. 3: Увеличение количества вложений в граф с ростом вершин в паттерне.

выполняться. При добавлении ребра в паттерн количество вложений не увеличится, следовательно, не увеличится и значение метрики. Если же добавляется новая вершина, то количество вложений может увеличиться, остаться прежним или уменьшиться. Последние два случая аналогичны варианту с добавлением ребра. Рассмотрим теперь случай, когда количество вложений увеличилось. Тогда результат метрики не сможет стать больше. Действительно, ведь предыдущее значение содержалось среди старых вершин, а после добавления новой количество их образов не сможет увеличиться, следовательно, как максимум значение метрики для нового паттерна будет принадлежать одной из них, что не больше, чем значение старого паттерна.

Заметим, что естественно определить эту метрику и для паттернов $\text{supp}(Q, G)$. Для этого можно пренебречь условием выполнимости зависимости. Это понадобится для более эффективной реализации алгоритма.

4.2. Предобработка

Прежде чем перейти к работе алгоритма, понадобится произвести некоторые подготовительные работы. Для этого необходимо обойти начальный граф, запомнить все метки вершин и рёбер, которые в нём присутствуют, все вложения одновершинных паттернов, а так же сгенерировать структуру, упрощающую генерацию литералов для паттер-

нов.

Информация о метках пригодится для генерации паттернов, ведь паттерны с иными метками совершенно точно не будут иметь вложений в данный граф.

По ходу просмотра будут сразу генерироваться паттерны, состоящие из одной вершины, а так же записываться все вложения для каждого из них. Полученные данные будут использоваться как вход на первую итерацию алгоритма.

Структура для генерации литералов представляет собой ассоциативный массив. Его ключами являются метки вершин, а значениями — новые ассоциативные массивы, отвечающие за возможные атрибуты этих вершин. Такое разделение сделано из соображений зависимости атрибутов от меток вершин. Во вложенных массивах в качестве ключей выступают названия атрибутов, а в качестве значений — списки значений этих атрибутов, которые встречаются в графе. Пример такой структуры представлен на Рис. 4.

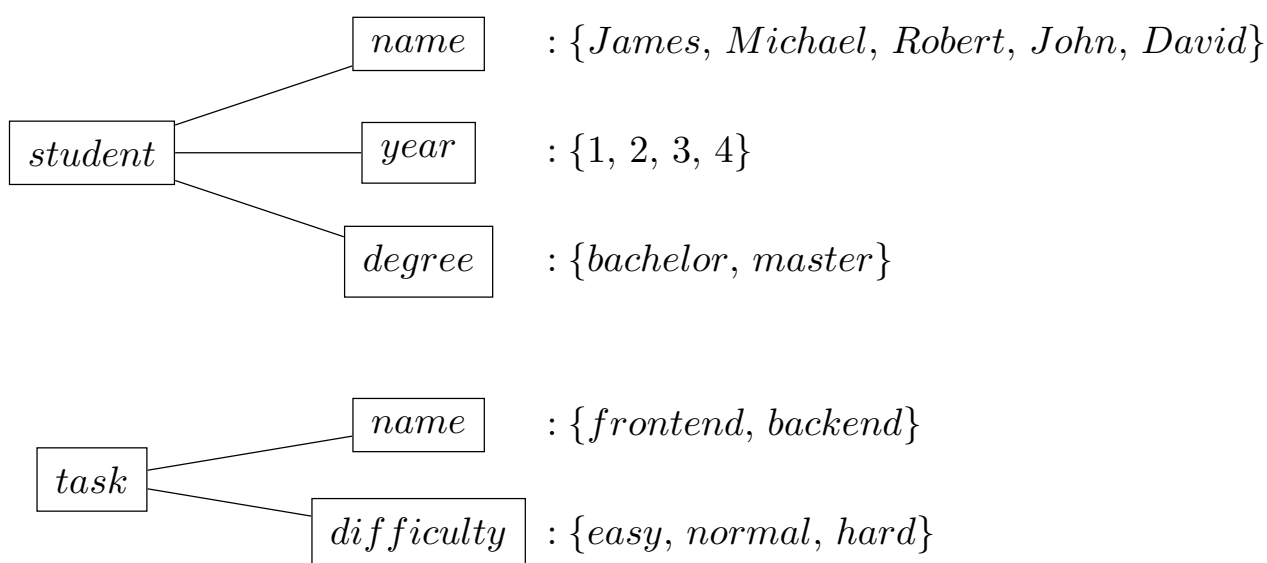


Рис. 4: Структура для генерации литералов.

Как видно из примера, названия атрибутов у разных вершин могут совпадать, однако списки значений у них могут отличаться. Учитывая эту информацию, генерация литералов будет происходить гораздо оптимальнее.

Обозначим за Θ множество всех вершинных меток ($L : V \rightarrow \Theta$), а за Ψ — множество всех рёберных меток ($L : E \rightarrow \Psi$).

Рассмотрим пример. Пусть, $\Theta = \{a, b, c\}$, $\Psi = \{e\}$. Тогда набор всех сгенерированных паттернов будет таким, как показано на Рис. 5.

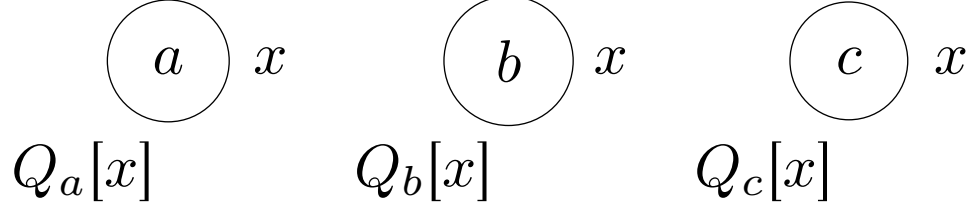


Рис. 5: Пример: начальные паттерны.

Далее вычисляется функция *supp* для каждого сгенерированного паттерна. Предположим, пользователь ввёл значение для $\sigma = 100$. А результаты применения функций оказались следующими:

- $\text{supp}(Q_a, G) = 216$.
- $\text{supp}(Q_b, G) = 97$.
- $\text{supp}(Q_c, G) = 54$.

Тогда только паттерн $Q_a[x]$ останется для первой итерации, так как $\text{supp}(Q_a, G) > \sigma = 100$, а остальные рассматриваться не будут. Это позволяет существенно сократить время работы алгоритма.

4.3. Основная работа

Алгоритм является итеративным. В целом, каждая итерация состоит из горизонтальной и вертикальной генерации графовых зависимостей. Она обновляет список паттернов, а так же списки вложений и запрещённых правил для каждого паттерна. Для начальных паттернов списки запрещённых литералов являются пустыми.

Горизонтальная генерация получает на вход список паттернов, создаёт всевозможные литералы для каждого из них с помощью структуры литералов, и проверяет, выполняются ли они, при помощи списка

вложений. Назовём этот набор литералов буквой Γ . Из множества Γ выделяется литерал-заключение l , а $\Gamma := \Gamma \setminus \{l\}$. Процедура генерации использует древовидную структуру, в узлах которой содержатся множества литералов. На самом верхнем уровне записывается единственное пустое множество. Потом происходит проверка зависимости с данным паттерном и правилом вида $(\emptyset \rightarrow l)$. Если оно выполнено, то работа окончена, переходим к следующему паттерну. Иначе достраиваем дерево: на втором, нижнем, уровне записываем все одноэлементные множества из Γ , ставим ребро из пустого множества к написанным. Проверка осуществляется с учётом параметра σ и параллельным подсчётом метрики *supp*. Если для какого-то паттерна нашлась GFD $\varphi = P[X \rightarrow Y]$ такая, что $\text{supp}(\varphi, G) > \sigma$ и $G \models \varphi$, φ добавляется в результат, а правило $X \rightarrow Y$ добавляется в список запрещённых правил для этого паттерна. Если нашли такой литерал, на котором GFD выполнялась, дальше вниз от этого узла прекращается построение дерева. В ином случае на третьем уровне пишутся двухэлементные множества, не содержащие выполненный литерал. Рёбра рисуются от узлов к узлам, подмножествами которых они являются. Процедура повторяется, пока нельзя будет сгенерировать следующий слой, или если дошли до множеств небольшой длины, например, 3-4. Такая эвристика обусловлена тем, что пользователю маловероятно, что понадобится сложные правила с большим количеством литералов, так как они почти наверняка не будут выполнены и весьма не информативны.

Рассмотрим на примере. Допустим, $\Gamma = \{l_1, l_2, l_3, l_4\}$, и l — заранее выбранный литерал, который будет содержаться в заключении правила.

Тогда дерево будет иметь вид, изображённый на Рис. 6. Где GFD $Q_a[x](l_1 \rightarrow l)$ и $Q_a[x](\{l_2, l_4\} \rightarrow l)$ удовлетворяют графу G .

Как видно, все генерируемые правила будут всегда содержать в правой части один литерал. После того, как все зависимости будут получены, их можно упростить, склеивая две, содержащие в левых частях одно и то же, пока такие не закончатся.

Далее происходит вертикальная генерация. Она заключается в со-

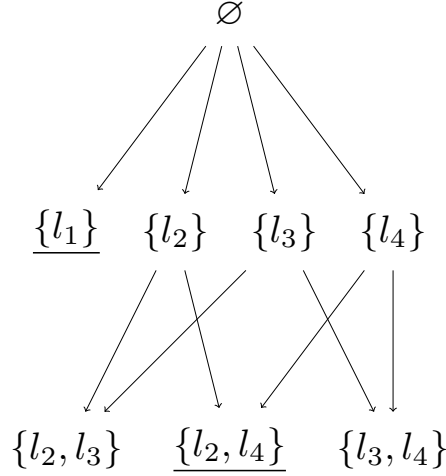


Рис. 6: Пример: дерево литералов.

здании паттернов для следующей итерации. Для каждого паттерна создаётся новый посредством добавления ребра. Ребро может быть добавлено с новой вершиной или без неё. Как только паттерн сгенерирован, происходит проверка на изоморфность с помощью, например, алгоритма ульмана [6] с новыми уже сгенерированными паттернами. Если такой нашёлся, то производится добавление к списку запрещённых правил найденного паттерна запрещённых правил сгенерированного. Иначе список запрещённых правил копируется от того, на основе которого был сгенерирован текущий, а так же создаётся новый список вложений на основе него.

После того, как все новые паттерны сгенерированы, они пропускаются через фильтр метрикой *supp*. Очевидно, что

$$supp(P[X \rightarrow Y], G) \leq supp(P, G),$$

Следовательно, если $supp(P, G) < \sigma$, то паттерн P можно дальше не рассматривать.

За счёт того, что когда находится графовая зависимость $\varphi = P[X \rightarrow Y]$, выполняющаяся на данном графе, все зависимости $\psi = Q[X \rightarrow Y]$, такие, что $P \subset Q$, запрещены, так как для них правило $X \rightarrow Y$ попадает в список запрещённых, это помогает искать лишь минимальные графовые зависимости.

Алгоритм завершается, как только был сгенерирован пустой список новых паттернов или после k -ой итерации.

5. Реализация

5.1. Алгоритм

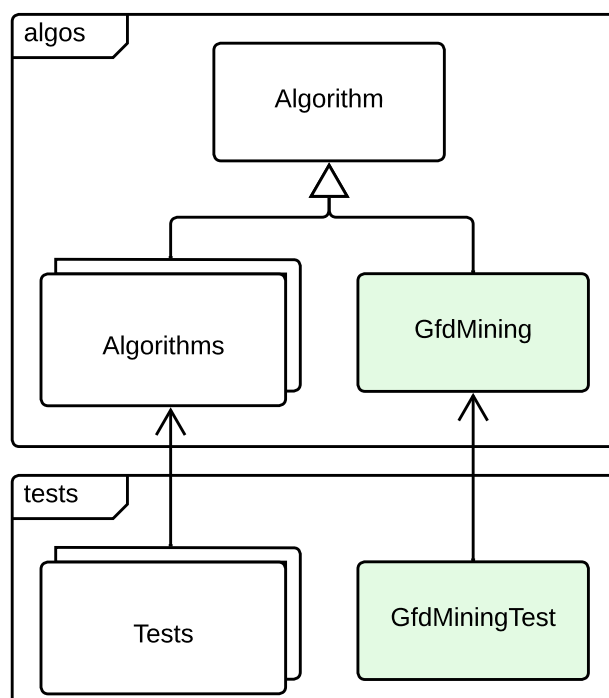


Рис. 7: Диаграмма классов для реализованного алгоритма.

Диаграмма на Рис. 7 показывает, как выглядит описанная реализация в проекте Desbordante. Все алгоритмы лежат в пространстве имён `algos` и отнаследованы от класса `Algorithm`, который описывает общий интерфейс для всех алгоритмов, работающих с функциональными зависимостями. Зелёным выделено то, что было реализовано.

При реализации алгоритма была использована графовая библиотека `boost`. Она позволяет гибко определить свой собственный тип графа, названный в проекте `graph_t`, его свойства, а также свойства вершин и рёбер. Так как графовые зависимости требуют от графа быть графом с атрибутами, все вершины были сконструированы содержать в себе список атрибутов. Помимо этого все элементы графа имеют собственные метки, а графы являются неориентированными.

В качестве литерала было удобно использовать структуру `std::pair`, так как литерал состоит из двух частей, к тому же такие пары легко

сравнивать между собой.

Графовые зависимости представлены в виде класса *Gfd*, которые содержат в себе граф типа *graph_t*, а также два множества: посылку и заключение, которые представляют собой списки литералов.

Для алгоритма были реализованы несколько вспомогательных функций:

- *GenerateLiterals()*, которая позволяет по входящему паттерну построить список литералов с помощью вышеописанной структуры литералов, которые возможны на данном паттерне.
- *Supp()*, позволяющая считать метрику частоты встречаемости для паттерна или зависимости, реализация которой была подробно описана в разделе “Входные параметры”.
- *AddEdge()*, которая генерирует все возможные паттерны по данному, имеющие хотя бы одно вложение в граф, путём добавления одного ребра.
- *AddVertex()* аналогична функции *AddEdge()*, за исключением того, что она создаёт новые паттерны через добавление новой вершины к паттерну.

Кроме них была реализована функция *Initialize()*, проводящая работу, описанную в разделе “Предобработка”.

Функция, выполняющая основную работу, состоит из большого цикла, который выполняется, пока множество сгенерированных на прошлой итерации паттернов не пусто, или достигнуто максимальное число итераций k . В теле цикла сначала вызывается функция *HorizontalSpawn()*, производящая правила литералов для каждого паттерна, генерирующая графовые зависимости, а затем проверяющая их выполнимость. После её вызова происходит генерация паттернов для следующей итерации с помощью функций *AddEdge()* и *AddVertex()*, и выполняется фильтрация тех паттернов, чьё количество вложений меньше входного параметра σ через вызов функции *FilterSupp()*.

5.2. Тестирование

Для апробации алгоритма было добавлено 5 тестов на небольших графах в среднем имеющих 13 вершин и 19 рёбер.

Тесты направлены на:

- Выявление минимальных графовых зависимостей.
- Обработку графов, не имеющих зависимостей.
- Получение графовых зависимостей с заключениями, имеющими более одного литерала.
- Корректную работу на синтетических данных.

5.3. Python bindings

Чтобы можно было выполнять C++ код алгоритма поиска функциональных зависимостей в графах, была использована библиотека `pybind11`². Это простая библиотека для преобразования типов данных между Python и C++. Она в основном применяется для создания привязок Python к уже существующему коду на C++.

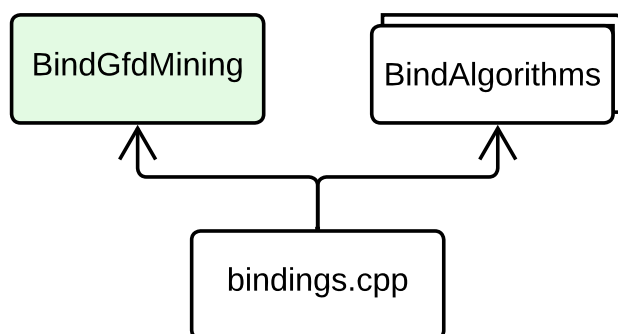


Рис. 8: Схема привязок алгоритмов Desbordante к Python.

На Рис. 8 показана схема привязки C++ кода к Python. Макрос в файле `bindings.cpp` позволяет определить Python-модуль `desbordante`. Также был создан новый файл `BindGfdMining`, предоставляющий API для работы с алгоритмом поиска графовых зависимостей на Python.

²<https://github.com/pybind/pybind11> (дата обращения 9.03.2024)

5.4. Примеры

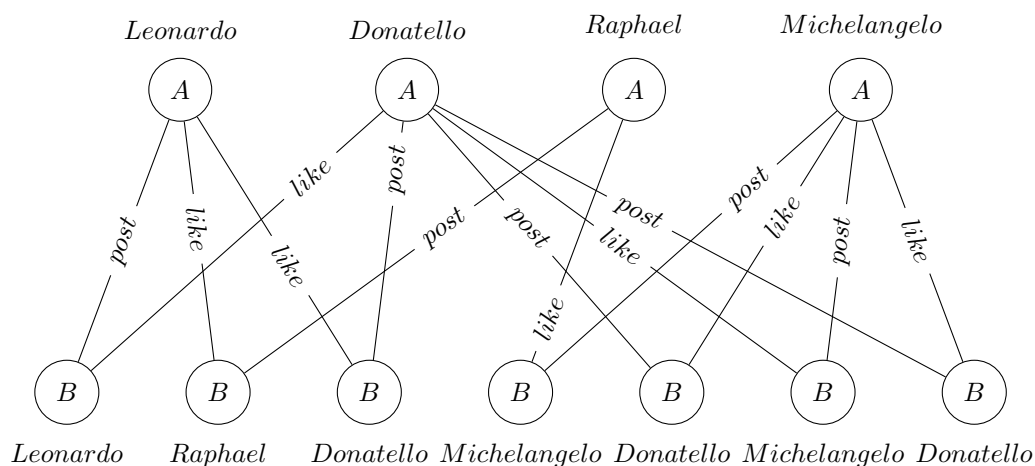


Рис. 9: Пример первого графа.

Для проверки работы привязок к Python были написаны следующие примеры. Рассмотрим первый. Граф изображён на Рис. 9. Здесь использовались следующие сокращения: A — *account*, B — *blog*. У вершин с меткой A есть атрибут *name*, показывающий никнейм; у вершин с меткой B — *author*, говорящий о том, кто написал данный блог. Рядом с вершинами подписаны значения этих атрибутов. Рёбра так же имеют метки: *post*, означающую то, кем был написан блог, и *like*, означающую одобрение другого человека.

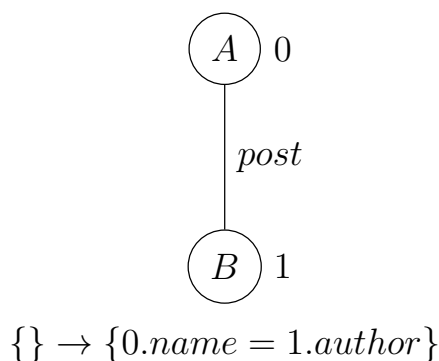


Рис. 10: Результат работы алгоритма на первом графе.

Найденная зависимость (Рис. 10) говорит о том, что если автор запостил блог, то в авторстве этого блога всегда указано имя запостившего его человека.

Рассмотрим теперь второй пример.

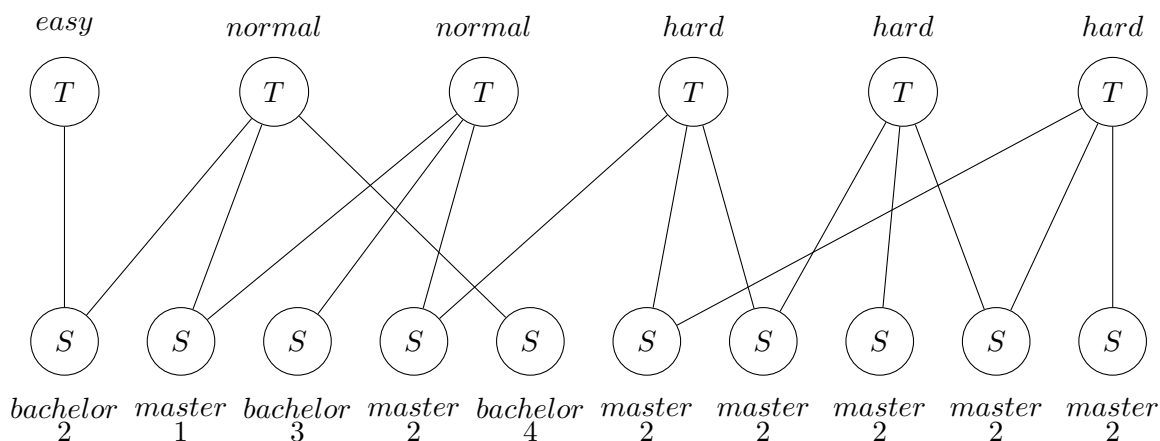
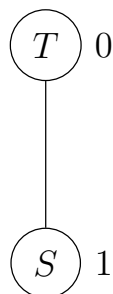


Рис. 11: Пример второго графа.

На Рис. 11 представлен граф с вершинами двух типов: T (*task*) и S (*student*). У вершин с меткой T есть два атрибута: *name* и *difficulty*. Они показывают название задачи и её сложность в общем понимании. У вершин типа S — три атрибута: *name*, *degree* и *year*. Они обозначают имя студента, уровень получаемого образования и курс. Значения всех атрибутов подписаны рядом с соответствующими вершинами, кроме атрибута *name*, так как этот атрибут не несёт значимой информации.



$$\{0.difficulty = hard\} \rightarrow \{1.degree = master \ \& \ 1.year = 2\}$$

Рис. 12: Результат работы алгоритма на втором графе.

Рассмотрим найденную графовую зависимость, представленную на Рис. 10. Она говорит о том, что над сложной задачей трудятся только магистранты второго курса.

Данные примеры были успешно протестированы с помощью Python, и были получены ожидаемые результаты.

5.5. Python CLI

Для взаимодействия с python-модулем `desbordante` используется модуль `Click`³. Это инструмент, который позволяет создавать интерфейсы командной строки.

`Click` выделяется среди аналогичных инструментов тем, что он предоставляет возможность легко и быстро создавать опции командной строки с использованием минимального количества кода. Кроме того, он предлагает инструменты для настройки и обеспечивает лёгкое масштабирование кода.

Командная строка `desbordante` имеет несколько основных опций, которые необходимы для запуска алгоритмов. Опция `help` предоставляет информацию о том, как использовать командную строку. Опции `task` и `algo` отвечают за задачу и конкретный алгоритм, который нужно запустить. Важно отметить, что одну и ту же задачу могут выполнять несколько алгоритмов. Например, для валидации графовых зависимостей доступны наивный, базовый и эффективный алгоритмы проверки графовых зависимостей [7].

Каждый алгоритм имеет свой уникальный набор опций. В случае алгоритма поиска графовых зависимостей этот набор включает путь к `dot`-файлу, содержащему представление графа, а также целые числа k и σ , которые являются параметрами алгоритма. Набор опций создаётся с помощью декоратора.

Пример вызова валидатора через Python консоль:

```
$ desbordante --task=gfd_mining --algo=gfd_miner  
--graph=examples/graph.dot --gfd_k=3 --gfd_sigma=10
```

³<https://click.palletsprojects.com/en/8.1.x> (дата обращения 24.03.2024)

Заключение

По итогам работы инструментарий Desbordante был расширен возможностью поиска графовых функциональных зависимостей в графовых структурах, а также появилась возможность запускать реализованный алгоритм на Python и через консоль.

Результаты работы:

- Реализован алгоритм майнинга графовых функциональных зависимостей и произведено тестирование.
- Реализована возможность запускать алгоритм поиска графовых зависимостей из скриптов, написанных на языке программирования Python.
- Созданы скрипты-примеры работы алгоритма поиска графовых зависимостей на языке программирования Python.
- Обеспечена возможность запускать алгоритм поиска графовых зависимостей через консоль путём реализации соответствующей подсистемы.

Исходный код алгоритма, Python-привязок и примеров доступен по ссылке⁴. Ссылка на код, обеспечивающий возможность запускать реализованный алгоритм через консоль⁵.

⁴<https://github.com/Desbordante/desbordante-core/pull/465> (дата обращения 29.09.2024)

⁵<https://github.com/Desbordante/desbordante-cli/pull/5> (дата обращения 02.10.2024)

Список литературы

- [1] Fan Wenfei, Hu Chunming, Liu Xueli, and Lu Ping. Discovering Graph Functional Dependencies. — 2020. — Access mode: <https://dl.acm.org/doi/abs/10.1145/3397198> (online; accessed: 2022-10-16).
- [2] Bi Fei, Chang Lijun, Lin Xuemin, Qin Lu, and Zhang Wenjie. Efficient Subgraph Matching by Postponing Cartesian Products. — 2016. — Access mode: <https://dl.acm.org/doi/abs/10.1145/2882903.2915236> (online; accessed: 2023-02-23).
- [3] Fan Wenfei, Liu Xueli, and Cao Yingjie. Parallel Reasoning of Graph Functional Dependencies. — 2018. — Access mode: <https://ieeexplore.ieee.org/abstract/document/8509281> (online; accessed: 2022-10-17).
- [4] Fan Wenfei, Wu Yinghui, and Xu Jingbo. Functional Dependencies for Graphs. — 2016. — Access mode: <https://dl.acm.org/doi/abs/10.1145/2882903.2915232> (online; accessed: 2022-09-14).
- [5] Chernikov Anton, Litvinov Yurii, Smirnov Kirill, and Chernishev George. FastGFDs: Efficient Validation of Graph Functional Dependencies with Desbordante. — 2023. — Access mode: <https://elibrary.ru/item.asp?id=53943942> (online; accessed: 2024-03-09).
- [6] Ullmann J. R. An Algorithm for Subgraph Isomorphism. — 1976. — Access mode: <https://dl.acm.org/doi/abs/10.1145/321921.321925> (online; accessed: 2022-11-11).
- [7] Черников Антон. Реализация эффективного алгоритма проверки графовых функциональных зависимостей в платформе Desbordante. — 2023. — Access mode: <https://github.com/Desbordante/desbordante-core/blob/main/docs/papers/FastGFDs%20-%20Anton%20Chernikov%20-%20BA%20thesis.pdf> (online; accessed: 2024-03-24).

- [8] Черников Антон. Расширение возможностей профилировщика данных Desbordante по работе с графовыми зависимостями. — 2024. — Access mode: <https://github.com/Desbordante/desbordante-core/blob/main/docs/papers/Bindings%2C%20examples%2C%20CLI%20and%20mining%20review%20for%20GFD%20-%20Anton%20Chernikov%20-%202023%20autumn.pdf> (online; accessed: 2024-10-03).