

# 1. Введение

В этом отчёте производится сравнение двух алгоритмов поиска частых наборов бинарных данных - AD и DepthProject. На вход подается булева матрица  $L = (a_{ij})$  размера  $m \times n$ , в которой элемент  $a_{ij}$  равен 1, если в транзакции с номером  $i$  значение порядка равно 1. Подается минимальная поддержка  $s$  - целое число в диапазоне от 1 до  $m$ .

## 2. Алгоритм AD

### 2.1. Описание

Алгоритм AD строит частые наборы пошагово. Введём обозначения:  $R(L)$  - множество всех столбцов матрицы  $L$ ,  $S_1(L)$  - множество столбцов из  $R(L)$  с не менее чем  $s$  единицами,  $e_1(R)$  и  $e_2(R)$  - столбцы соответственно с наименьшим и наибольшим номерами в  $R \subseteq R(L)$ . Для набора  $Q_t = j_1, \dots, j_t$  обозначим  $R_t$  - множество столбцов, номера которых больше  $j_t$ ,  $G_t$  - множество столбцов  $h$  из  $R_t$ , которые в объединении с  $Q_t$  дают совместимый набор.

Сначала строится первый частый набор  $Q = e_1(S_1(L))$ . Затем для каждого  $Q = \{j_1, \dots, j_r\} \neq \{e_2(S_1(L))\}$  строится следующий за ним набор  $\Delta Q$  по следующим правилам:

- 1)  $G_r \neq \emptyset$ :  $\Delta Q = Q \cup \{e_1(G_r)\}$ ;
- 2)  $G_r = \emptyset$ :
  - а)  $r = 1$ :  $\Delta Q = \{e_1(S_r)\}$ ;
  - б)  $r > 1$  и  $G_{r-1} \neq \emptyset$ :  $\Delta Q = Q_{r-1} \cup e_1(G_{r-1})$ ;
  - в)  $r > 1$  и  $G_{r-1} = \emptyset$ : если  $r = 2$ , то  $\Delta Q = \{e_1(R_r)\}$ , иначе  $\Delta Q = Q_{r-2} \cup \{e_1(R_{r-1})\}$ ;

Таким образом, алгоритм обходит дерево решений, вершины которого - частые наборы. Максимальные наборы находятся среди висячих вершин. В терминах алгоритма это значит, что если  $G_t = \emptyset$ , то  $Q_t$  является либо максимальным частым набором, либо подмножеством ранее найденного максимального частого набора. Проверку на максимальность можно проводить как по ранее найденным наборам (тогда сложность алгоритма зависит от выхода задачи), так и по базе данных (модификация, предложенная Н. А. Драгуновым). Второй способ значительно быстрее.

Частый наборы выводятся в алфавитном порядке.

### 2.2. Реализация

Сначала строятся первый частый набор  $Q[1] = \{e_1(R(L))\}$  и последний частый набор  $L = \{e_2(R(L))\}$ . Затем последовательно строятся наборы  $Q[2], Q[3], \dots$  из одной ветви дерева решений вплоть до первого максимального набора. После его нахождения начинает работать основной цикл, выход из которого осуществляется если очередной найденный набор равен  $L$ .

Используется процедура Обновить( $T, Q$ ), которая:

- 1) Зануляет строки матрицы  $T$ , дающие в пересечении со столбцами из  $Q$  нули.
- 2) Удаляет столбцы с номерами из  $Q$ .

Таким образом, первый совместимый столбец модифицированной матрицы есть первое частое расширение набора  $Q$ .

Ниже представлен псевдокод алгоритма AD:

---

**Алгоритм 1** DepthProject

---

```

1: Процедура AD(БазаДанных :  $T$ , Поддержка :  $s$ )
2:    $Q = \emptyset$ 
3:    $L = \text{ПоследнийЧастыйСтолбец}(T, s)$ 
4:    $t = \text{ПервыйЧастыйСтолбец}(T, s)$ 
5:   До тех пока  $t \neq \emptyset$  выполнять
6:      $Q = Q \cup \{e_1(G_r)\}$   $\triangleright r = |Q|$ 
7:      $Q$  — частый набор
8:      $t = \{e_1(G_r)\}$ 
9:   Конец цикла
10:   $T' = T$ 
11:  До тех пока  $Q \neq L$  выполнять
12:     $T' = \text{Обновить}(T', Q)$ 
13:     $t = \text{ПервыйЧастыйСтолбец}(T', s)$ 
14:    Если  $t \neq \emptyset$  тогда
15:       $Q = Q \cup \{t\}$ 
16:    иначе
17:       $Q$  — частый набор
18:      Проверка  $Q$  на максимальность
19:      Если  $|Q| = 1$  тогда
20:         $Q = \text{ПервыйЧастыйСтолбец}(T', Q_1, s)$ 
21:      иначе
22:        Если  $|Q| = 2$  тогда
23:           $Q = \text{ПервыйЧастыйСтолбец}(T', Q_1, s)$ 
24:        иначе
25:           $p = \text{ПредпоследнийЭлемент}(Q)$ 
26:           $T' = \text{Обновить}(T, Q_{r-2})$ 
27:           $t = \text{ПервыйСовместимыйСтолбец}(T', p, s)$ 
28:           $Q = Q \cup \{t\}$ 
29:        Конец условия
30:      Конец условия
31:    Конец условия
32:  Конец цикла
33: Конец процедуры

```

---

## 3. Алгоритм DepthProject

### 3.1. Пару слов о дереве решений

В DepthProject производится рекурсивный обход дерева решения. Введём некоторые обозначения. Пусть имеется дерево решений  $D_L$ :

1) Каждый узел соответствует частому набору. Корень дерева соответствует пустому набору элементов и обозначается  $Null$ .

2) Узел  $P = \{i_1, \dots, i_k\}$  имеет предка  $Q = \{i_1, \dots, i_{k-1}\}$ .

Обозначим за  $E(P)$  множество частых расширений узла  $P$ . Обозначим за  $F(P)$  множество частых расширений родительского узла  $Q$ , номера которых превышают последний элемент в  $P$ . Справедливы вложения:  $E(P) \subseteq F(P) \subset E(Q)$ .

Алгоритм обходит это дерево решений и выводит наборы в алфавитном порядке. Древовидная структура позволяет на каждом шаге сужать область поиска частых расширений, рассматривая так называемых кандидатов. Преимущество такого подхода в том, что кандидаты ищутся за константное время по предыдущему узлу.

### 3.2. Описание

Вызывается процедура DepthFirst, на вход которой подаются узел  $N$ , база данных  $T$  и маска  $B$ . При первом запуске  $N = Null$ ,  $T$  - вся база данных (т.е. вся бинарная матрица), маска  $B = (1, 1, \dots, 1)$  - двоичный вектор длины  $m$ . В каноническом изложении процедура выглядит следующим образом:

---

#### Алгоритм 2 DepthProject

---

```
1: Процедура DEPTHFIRST(Узел :  $N$ , БазаДанных :  $T$ , Маска :  $B$ )
2:    $C = F(P)$                                 ▷ Множество кандидатов на частое расширение
3:    $E = E(P)$                                 ▷ Множество частых расширений  $E = \{i_1, \dots, i_{|E|}\}$ 
4:   Положить  $N \cup i_r$  в узлы ДР
5:    $B' = \text{ПостроитьМаску}(N, B, T)$ 
6:   Если Условие проекции тогда
7:      $T' = \text{Проекция}(T, E, N, B')$ 
8:   иначе
9:      $T' = T$ 
10:  Конец условия
11:  Цикл <г := 1 до  $|E|$ > выполнять
12:     $\text{DepthFirst}(N \cup \{i_r\}, T', B')$ 
13:  Конец цикла
14: Конец процедуры
```

---

$F(P)$  - это номера столбцов из  $E(Q)$ , которые имеют номер, больший последнего номера в узле  $N$ . Здесь  $Q$  - ближайший предок вершины  $P$ . Маска  $B$  имеет 1 в  $i$ -й позиции тогда и только тогда, когда все столбцы текущего узла имеют единицу в  $i$ -й позиции. Проекция позволяет снизить время работы алгоритма. Вообще говоря,

это аналог функции "Обновить" в реализации алгоритма AD с той лишь разницей, что строки не зануляются, а удаляются. Уменьшая на каждом шаге базу данных  $T$  мы сужаем область поиска частых расширений  $E(P)$ . Также важно построить хороший метод поиска частых расширений.

### 3.3. Реализация

В текущей реализации DepthProject используется древовидная структура. Первым шагом находится множество кандидатов  $F(P)$  - эта операция производится быстро и сокращает дальнейший поиск частых расширений  $E(P)$ . Маска готовится перед вызовом этих функций. Здесь маска - это логическое умножение столбцов текущего узла. Для её нахождения не обязательно перемножать все столбцы, достаточно на каждом шаге умножать текущую маску на последний столбец в узле. Итак, алгоритм примет вид:

---

#### Алгоритм 3 DepthProject (реализованный)

---

```

1: Процедура DEPTHFIRST(Узел :  $N$ , Маска :  $B$ )
2:    $B = \text{ПостроитьМаску}(N, B, T)$ 
3:    $C = E(N)$ 
4:    $E = E(N, T, B, C)$ 
5:   Цикл  $\langle r := 1 \text{ до } |E| \rangle$  выполнять
6:      $N \cup \{i_r\} \rightarrow$  множество частых наборов
7:      $\text{DepthFirst}(N \cup \{i_r\}, T', B')$ 
8:   Конец цикла
9:   Если тогда  $E = \emptyset$ 
10:    Если тогда  $N$  — максимальный набор
11:     $N \rightarrow$  множество максимальных наборов
12:  Конец условия
13: Конец условия
14: Конец процедуры

```

---

В этой реализации нет проекции, ведь можно обойтись без неё и даже получить выигрыш в скорости. В самом деле, зная маску и последний элемент в текущем узле мы легко найдём следующий совместимый столбец по исходной базе данных и потратим столько же времени, сколько занимает одна проекция. Возможно, есть способ сделать проекцию быстрее поиска с маской, но преимущество в скорости должно быть достаточно большим, поскольку сам факт рекурсивной передачи столь большого массива данных занимает определённое время. На практике одна лишь передача базы данных без последующей проекции увеличивает время работы алгоритма в два раза. Так или иначе, ни на количество шагов, ни на сложность шага это не влияет. Проекция - программистское ухищрение авторов алгоритма, которое позволило им получить прирост в скорости. В моей реализации преимущество в скорости даёт маска и знание последнего элемента текущего узла.

## 4. Сравнение двух алгоритмов

При реализации алгоритмов использовалась одна и та же модель бинарной матрицы. В обоих алгоритмах маска - это бинарный вектор `bitset` (C++ 11 версии и выше).

На вход подаётся бинарная матрица,  $m$  - кол-во строк (транзакций),  $n$  - кол-во столбцов (элементов). Минимальная поддержка  $s = 0.1$ .

	<b>n = 10</b>	<b>n = 20</b>	<b>n = 30</b>
<b>m = 10</b>	183, 5 0.73 2.56	49557, 10 49 52	1029918, 10 1064 4250
<b>m = 100</b>	169, 76 0.89 1.52	1560, 785 3.39 6.52	9165, 4962 16.1 17.3
<b>m = 1000</b>	172, 117 1.34 1.87	1350, 1136 6.2 8.6	4518, 4053 23.3 33.9

В первой строке - кол-во частых наборов и кол-во макс. частых наборов. Во второй строке время работы AD в миллисекундах, в третьей - время работы DepthProject. Результаты были усреднены по 20 запускам.

Ниже приведена аналогичная таблица для минимальной поддержки  $s = 0.3$ .

	<b>n = 10</b>	<b>n = 20</b>	<b>n = 30</b>
<b>m = 10</b>	64, 10 0.184 0.398	573, 30 1.62 1.73	1225, 43 2.6 11.2
<b>m = 100</b>	16, 10 0.084 0.2	50, 34 0.383 0.216	110, 86 0.924 1.26
<b>m = 1000</b>	10, 10 0.154 0.222	20, 20 0.533 0.693	30, 30 0.817 1.043

Алгоритм AD немного быстрее, особенно в тех моментах, когда количество всех частых наборов много больше максимальных частых наборов. Нетрудно догадаться, что в этом случае дерево решений имеет большую глубину. Следовательно, рекурсивный вызов процедуры DepthProject имеет большую вложенность, а передача объектов в процедуру нагружает процессор. В этом и заключается преимущество последовательного подхода.

Изобразим на графике зависимость времени работы двух алгоритмов от количества столбцов при фиксированном количестве строк.

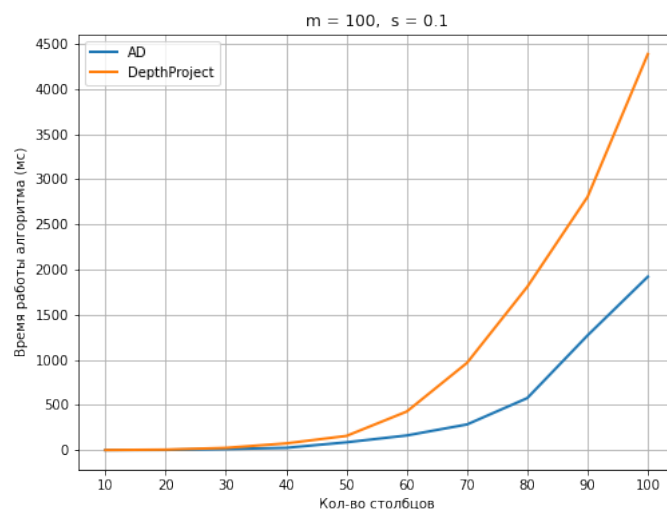


Рис. 1: Зависимость времени работы от кол-ва столбцов

Как видим, скорости счёта двух алгоритмов отличаются в константное число раз. Изобразим на графике отношение времени работы алгоритмов:

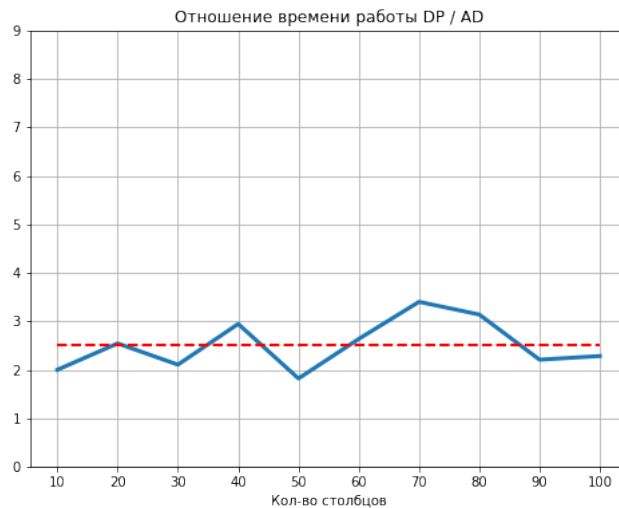


Рис. 2: Отношение времени работы DepthProject и AD

В среднем скорость счёта отличается в 2.5 раза. Это подтверждает тот факт, что алгоритмы совершают одинаковое число шагов. Сложность шага в DepthProject чуть больше, но это связано исключительно с реализацией.

Теперь изобразим на графике зависимость времени работы двух алгоритмов от количества строк при фиксированном количестве столбцов. Положим  $n = 20$ ,  $s = 0.1$ :

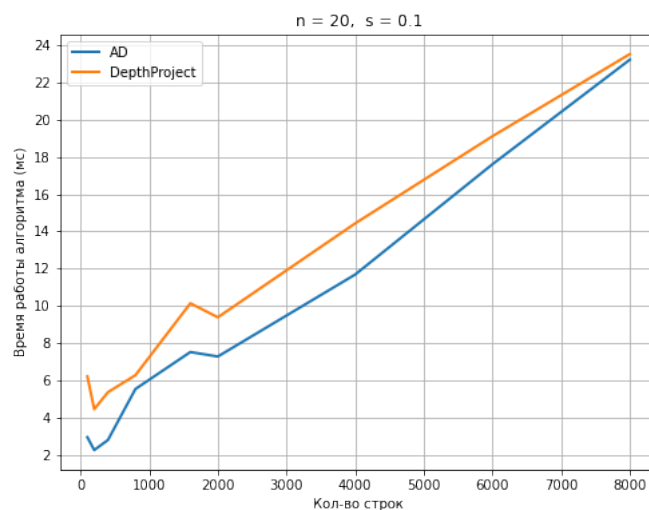


Рис. 3: Зависимость времени работы от кол-ва строк

Время работы двух алгоритмов асимптотически совпадает. Ниже представлен график отношения скоростей счёта:

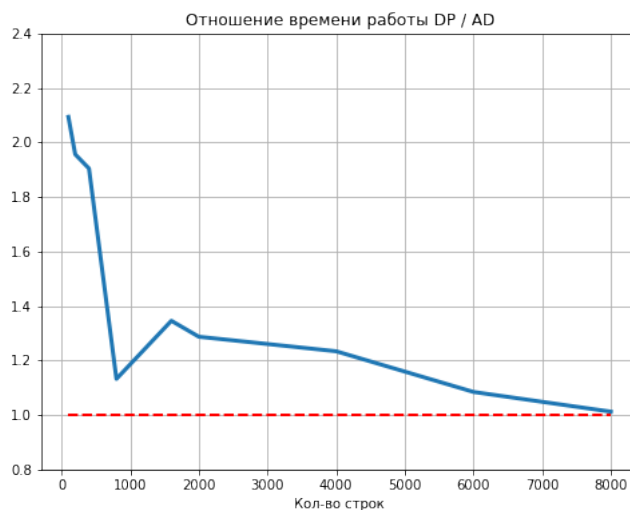


Рис. 4: Отношение времени работы DepthProject и AD

## 5. Вывод

Алгоритмы AD и DepthProject совершают одинаковое число шагов, имеют схожую сложность шага и выводят наборы в алфавитном порядке. Основное отличие заключается в том, что AD выполняется последовательно и использует меньше вычислительных ресурсов, в то время как DepthProject работает рекурсивно и строит полное дерево решений.

## 6. Дополнение

### 6.1. Удаление нечастых столбцов

Пусть дана бинарная матрица  $T \in B^{m \times n}$ ,  $B = \{0, 1\}$  и минимальная поддержка  $s \in [0, 1]$ . Мы можем удалить все столбцы матрицы  $T$ , имеющие меньше  $m \cdot s$  единиц, ведь они не войдут ни в один частый набор. Назовём такие столбцы нечастыми.

Обозначим среднее количество нечастых столбцов за  $Y$ . Пусть матрица  $T$  является случайной и берётся из равномерного распределения. Ясно, что при малых  $s$  нечастых столбцов почти не будет, а при  $s \approx 0.5$  и  $s > 0.5$  такие столбцы почти наверняка найдутся. Посчитаем, при каких  $s, m, n$  среднее количество таких столбцов отлично от нуля.

Элементы матрицы  $T$  имеют распределение Бернулли с  $p = 0.5$ . Следовательно, распределение единиц в отдельно взятом столбце представляет собой биномиальное распределение. Обозначим кол-во единиц в столбце за  $k$ , тогда

$$P(k = i) = C_m^i \cdot \left(\frac{1}{2}\right)^i \cdot \left(\frac{1}{2}\right)^{m-i} = \frac{1}{2^m} \cdot C_m^i.$$

Тогда несложно посчитать вероятность того, что столбец нечастый:

$$P(k < s) = \sum_{i=0}^{s-1} P(k = i) = \frac{1}{2^m} \cdot \sum_{i=0}^{s-1} C_m^i$$

Чтобы найти среднее количество нечастых столбцов умножим полученную вероятность на  $n$ . Окончательно получим:

$$Y = \frac{n}{2^m} \cdot \sum_{i=0}^{s-1} C_m^i$$

Можно показать, что  $Y > 1$  начиная с  $s > 0.4 \cdot m$  (примерно). Итак, если минимальная поддержка  $s > 0.4$ , то есть смысл перед поиском частых наборов удалить нечастые столбцы.