

Функциональное программирование

Ричард Бёрд

Жемчужины проектирования алгоритмов

Функциональный подход

С примерами на языке Haskell

УДК 004.021+004.421
ББК 32.973-018
Б11

Б11 Ричард Бёрд

Жемчужины проектирования алгоритмов: функциональный подход / Пер. с англ. В. Н. Брагилевского и А. М. Пеленицына. – М.: ДМК Пресс, 2013. – 330 с.: ил.

ISBN 978-5-94074-867-0

В этой книге Ричард Бёрд представляет принципиально новый подход к проектированию алгоритмов, а именно проектирование посредством формального вывода. Основное содержание книги разделено на 30 коротких глав, называемых жемчужинами, в каждой из которых решается конкретная программистская задача. Эти задачи, некоторые из них абсолютно новые, происходят из таких разнообразных источников, как игры и головоломки, захватывающие комбинаторные построения и более традиционные алгоритмы сжатия данных и сопоставления строк.

Каждая жемчужина начинается с постановки задачи, формулируемой на функциональном языке программирования Haskell, чрезвычайно мощном и в то же время лаконичном, позволяющем легко и просто выражать алгоритмические идеи. Новшество книги состоит в том, что каждое решение формально вычисляется из исходной постановки задачи посредством обращения к законам функционального программирования.

Издание предназначено для программистов, увлекающихся функциональным программированием, студентов, аспирантов и преподавателей, интересующихся принципами проектирования алгоритмов, а также всех, кто желает приобрести и развить навыки рассуждений в эквациональном стиле применительно к программам и алгоритмам.

УДК 004.021+004.421
ББК 32.973-018

Все права защищены. Любая часть этой книги не может быть воспроизведена в какой бы то ни было форме и какими бы то ни было средствами без письменного разрешения владельцев авторских прав.

Материал, изложенный в данной книге, многократно проверен. Но, поскольку вероятность технических ошибок всё равно существует, издательство не может гарантировать абсолютную точность и правильность приводимых сведений. В связи с этим издательство не несёт ответственности за возможные ошибки, связанные с использованием книги.

ISBN 978-0-521-51338-8 (англ.)

ISBN 978-5-94074-867-0 (рус.)

© 2010 Cambridge University Press

© Перевод на русский язык, оформление,
ДМК Пресс, 2013

Оглавление

<i>Предисловие</i>	9
1 Наименьшее отсутствующее число	12
2 Превосходная задача	19
3 Улучшаем седловой поиск	24
4 Задача о выборке	35
5 Сортировка попарных сумм	42
6 Делаем сотню	49
7 Строим дерево минимальной высоты	58
8 Распутываем жадные алгоритмы	68
9 Поиск знаменитостей	75
10 Удаляем повторы	84
11 Вовсе не максимальная сумма сегмента	94
12 Ранжируем суффиксы	101
13 Преобразование Барроуза–Уилера	115
14 Последний хвост	128
15 Все общие префиксы	139

16	Алгоритм Бойера—Мура	145
17	Алгоритм Кнута—Морриса—Пратта	156
18	Планирование в «Час пик»	166
19	Простой алгоритм решения sudoku	178
20	Задача «Обратного отсчёта»	189
21	Хиломорфизмы и нексусы	202
22	Три способа вычисления определителей	215
23	Внутри выпуклой оболочки	224
24	Рациональное арифметическое кодирование	236
25	Целочисленное арифметическое кодирование	247
26	Алгоритм Шора—Вейта	262
27	Упорядоченная вставка	274
28	Бесцикловые функциональные алгоритмы	287
29	Алгоритм Джонсона—Троттера	298
30	Прядение паутины для чайников	306
	<i>Предметный указатель</i>	326

1

Наименьшее отсутствующее число

Введение

Рассмотрим задачу отыскания наименьшего натурального числа, отсутствующего в заданном конечном множестве натуральных чисел X . Здесь мы имеем дело с упрощённой версией более общей программистской задачи, в которой числа соответствуют некоторым объектам, а X — множество объектов, используемых в настоящий момент. Задача заключается в том, чтобы найти некоторый неиспользуемый объект, например, с наименьшим именем.

Разумеется, решение задачи зависит от способа представления множества X . Если X задано списком без повторений, где элементы упорядочены в порядке возрастания, то решение очевидно: в последовательности элементов следует искать первый пропуск. Предположим, однако, что множество X задано списком различных чисел в произвольном порядке, например:

[08, 23, 09, 00, 12, 11, 01, 10, 13, 07, 41, 04, 14, 21, 05, 17, 03, 19]

Как бы вы стали искать наименьшее число, отсутствующее в этом списке?

Не сразу становится очевидным, что имеется решение линейной сложности, ведь невозможно выполнить сортировку произвольного числового списка за линейное время. Тем не менее, такое решение существует, и целью этой жемчужины будет описание двух возможных стратегий: одна из них основана на использовании массивов языка Haskell, а вторая на методе «разделяй и властвуй».

Решение с использованием массива

Определим спецификацию задачи с помощью функции *minfree*:

$$\begin{aligned} \text{minfree} &:: [\text{Nat}] \rightarrow \text{Nat} \\ \text{minfree } xs &= \text{head} ([0..] \setminus xs) \end{aligned}$$

Выражение $us \setminus vs$ обозначает список тех элементов us , которые остаются после удаления всех элементов vs :

$$\begin{aligned} (\setminus) &:: \text{Eq } a \Rightarrow [a] \rightarrow [a] \rightarrow [a] \\ us \setminus vs &= \text{filter} (\not\in vs) us \end{aligned}$$

Хотя функция *minfree* и работает, для списка длины n она требует в худшем случае $\Theta(n^2)$ операций. К примеру, чтобы получить результат вызова $\text{minfree} [n - 1, n - 2..0]$ понадобится вычислить $i \notin [n - 1, n - 2..0]$, где $0 \leq i \leq n$, поэтому в конечном итоге получится $n(n + 1)/2$ проверок на равенство.

Ключевой факт для обеих стратегий решения, с массивом и по методу «разделяй и властвуй», заключается в том, что в списке xs содержатся не все элементы из промежутка $[0.. \text{length } xs]$. Поэтому наименьшее число, отсутствующее в xs , является одновременно наименьшим числом, отсутствующим в $\text{filter} (\leq n) xs$, где $n = \text{length } xs$. Решение, основанное на массиве, использует этот факт для построения контрольного массива чисел, находящихся в $\text{filter} (\leq n) xs$. Контрольный массив с индексами от 0 до n содержит $n + 1$ логическое значение, инициализированное *False*. Для каждого элемента x из xs , такого, что $x \leq n$, элемент контрольного массива в позиции x устанавливается равным *True*. Наименьшее отсутствующее число находится после этого по первой же позиции, равной *False*. Таким образом, $\text{minfree} = \text{search} \cdot \text{checklist}$, где

$$\begin{aligned} \text{search} &:: \text{Array Int Bool} \rightarrow \text{Int} \\ \text{search} &= \text{length} \cdot \text{takeWhile id} \cdot \text{elems} \end{aligned}$$

Функция *search* принимает на вход массив логических значений, преобразует массив в список и возвращает длину наибольшего начального сегмента, содержащего только истинные элементы. Полученная длина и будет индексом первого вхождения значения *False*.

Одним из возможных вариантов определения функции *checklist* с линейной сложностью является использование функции *accumArray* из модуля *Data.Array* стандартной библиотеки языка Haskell. Тип этой функции несколько пугающий:

$$Ix\ i \Rightarrow (e \rightarrow v \rightarrow e) \rightarrow e \rightarrow (i, i) \rightarrow [(i, v)] \rightarrow Array\ i\ e$$

Ограничение $Ix\ i$ требует, чтобы типовая переменная i принадлежала классу *Index*, например, *Int* или *Char*, она используется для обозначения индексов или позиций в массиве. Первый аргумент это «суммирующая» функция, она перевычисляет элемент массива (типа e) на основе некоторого значения (типа v). Второй аргумент определяет начальное значение для элемента массива в очередной позиции. Третий аргумент это пара из начального и конечного индексов. Наконец, четвёртый аргумент это ассоциативный список пар индекс–значение. Функция *accumArray* строит массив, обрабатывая ассоциативный список слева направо, при этом она меняет определяемые очередным индексом элементы массива, присваивая им результаты вызова суммирующей функции для прежнего элемента и очередного значения из ассоциативного списка. Этот процесс выполняется за линейное по длине ассоциативного списка время в предположении, что суммированию достаточно константного.

Теперь функцию *checklist* можно реализовать как вызов функции *accumArray*:

```
checklist    :: [Int] → Array Int Bool
checklist xs = accumArray (∨) False (0, n)
              (zip (filter (≤ n) xs) (repeat True))
              where n = length xs
```

Эта реализация не требует, чтобы элементы списка xs не повторялись. Единственное ограничение в том, чтобы они были натуральными числами.

Интересно, что функцией *accumArray* можно воспользоваться для сортировки числового списка за линейное время при условии, что все его элементы принадлежат ограниченному диапазону $(0, n)$. Заменяем *checklist* на *countlist*, где

```
countlist    :: [Int] → Array Int Int
countlist xs = accumArray (+) (0, n) (zip xs (repeat 1))
```

Теперь $sort\ xs = concat\ [replicate\ k\ x\ | (x, k) \leftarrow countlist\ xs]$. В сущности, применяя *countlist* вместо *checklist*, функцию *minfree* можно определить как позицию первого вхождения нулевого элемента.

Приведённая выше реализация строит массив за один проход, пользуясь умной библиотечной функцией. Более прозаичный способ реализовать функцию *checklist* состоит в явной отметке всех вхождений шаг за шагом с

применением операции присваивания с константной сложностью. В языке Haskell это возможно, только если обработка массива выполняется в подходящей монаде, например, в монаде с состоянием. В следующем примере реализация *checklist* использует модуль *Data.Array.ST*:

```
checklist xs =
  runSTArray (do
    { a ← newArray (0, n) False;
      sequence [writeArray a x True | x ← xs, x ≤ n];
      return a })
  where n = length xs
```

Впрочем, такое решение не удовлетворит чисто функционального программиста, поскольку оно эксплуатирует традиционную императивную парадигму, хотя и в функциональных одеждах.

Решение по методу «разделяй и властвуй»

Обратимся теперь к стратегии «разделяй и властвуй». Идея в том, чтобы выразить $minfree (xs ++ ys)$ в терминах $minfree xs$ и $minfree ys$. Запишем некоторые свойства операции $\setminus\setminus$:

$$\begin{aligned} (as ++ bs) \setminus\setminus cs &= (as \setminus\setminus cs) ++ (bs \setminus\setminus cs) \\ as \setminus\setminus (bs ++ cs) &= (as \setminus\setminus bs) \setminus\setminus cs \\ (as \setminus\setminus bs) \setminus\setminus cs &= (as \setminus\setminus cs) \setminus\setminus bs \end{aligned}$$

Эти свойства аналогичны соответствующим свойствам теоретико-множественных операций, в которых объединение множеств \cup заменяется на $++$, а разность множеств \setminus на $\setminus\setminus$. Предположим теперь, что as и vs не пересекаются, т.е. $as \setminus\setminus vs = as$, и что bs и vs также не пересекаются, т.е. $bs \setminus\setminus vs = bs$. Из указанных свойств операций $++$ и $\setminus\setminus$ следует, что

$$(as ++ bs) \setminus\setminus (us ++ vs) = (as \setminus\setminus us) ++ (bs \setminus\setminus vs)$$

Выберем теперь некоторое натуральное число b и положим $as = [0..b-1]$ и $bs = [b..]$. Пусть, далее, $us = filter (<b) xs$ и $vs = filter (\geq b) xs$. Тогда as и vs окажутся непересекающимися, а значит, такими же будут bs и us . Следовательно,

$$\begin{aligned} [0..] \setminus\setminus xs &= ([0..b-1] \setminus\setminus us) ++ ([b..] \setminus\setminus vs) \\ &\mathbf{where} (us, vs) = partition (<b) xs \end{aligned}$$

Haskell обеспечивает эффективную реализацию функции *partition*, которая разбивает список на те элементы, которые удовлетворяют предикату *p*, и те, которые ему не удовлетворяют. Поскольку

$$\text{head } (xs \text{ ++ } ys) = \text{if } \text{null } xs \text{ then head } ys \text{ else head } xs$$

получаем, по-прежнему для любого натурального числа *b*, что

$$\begin{aligned} \text{minfree } xs &= \text{if } \text{null } ([0..b-1] \setminus us) \\ &\quad \text{then head } ([b..] \setminus vs) \\ &\quad \text{else head } ([0..] \setminus us) \\ &\text{where } (us, vs) = \text{partition } (<b) \text{ } xs \end{aligned}$$

Следующий вопрос: можно ли реализовать проверку $\text{null } ([0..b-1] \setminus us)$ более эффективно, чем прямым вычислением, которое требует квадратичного времени по длине *us*? Да, так как на входе список неповторяющихся натуральных чисел, таковым же является *us*, причём каждый элемент *us* меньше *b*. Поэтому

$$\text{null } ([0..b-1] \setminus us) \equiv \text{length } us == b$$

Заметим, что предыдущее решение не зависело от предположения, что данный список не содержит дубликатов, однако оно оказывается ключевым для эффективной реализации по методу «разделяй и властвуй».

Дальнейшее изучение кода *minfree* подсказывает, что следует обобщить *minfree* до функции, скажем, *minfrom*, которая определена так:

$$\begin{aligned} \text{minfrom} &:: \text{Nat} \rightarrow [\text{Nat}] \rightarrow \text{Nat} \\ \text{minfrom } a \text{ } xs &= \text{head } ([a..] \setminus xs) \end{aligned}$$

где предполагается, что каждый элемент *x* больше или равен *a*. Тогда, при условии, что *b* выбрано таким образом, чтобы длины *us* и *vs* были меньше длины *xs*, следующее рекурсивное определение *minfree* оказывается вполне обоснованным:

$$\begin{aligned} \text{minfree } xs &= \text{minfrom } 0 \text{ } xs \\ \text{minfrom } a \text{ } xs &\mid \text{null } xs &&= a \\ &\mid \text{length } us == b - a &&= \text{minfrom } b \text{ } vs \\ &\mid \text{otherwise} &&= \text{minfrom } a \text{ } us \\ &\text{where } (us, vs) = \text{partition } (<b) \text{ } xs \end{aligned}$$

Остаётся выбрать *b*. Ясно, что нам нужно $b > a$. К тому же, хотелось бы иметь *b* таким, чтобы максимум из длин *us* и *vs* был настолько малым,

насколько это возможно. Правильный выбор b , удовлетворяющего указанным требованиям, таков:

$$b = a + 1 + n \operatorname{div} 2$$

где $n = \operatorname{length} xs$. Если $n \neq 0$ и $\operatorname{length} us < b - a$, то

$$\operatorname{length} us \leq n \operatorname{div} 2 < n$$

И, если $\operatorname{length} us = b - a$, то

$$\operatorname{length} vs = n - n \operatorname{div} 2 - 1 \leq n \operatorname{div} 2$$

При таком выборе параметра b число операций $T(n)$, необходимое для вычисления $\operatorname{minfrom} 0 xs$, где $n = \operatorname{length} xs$, удовлетворяет рекуррентному соотношению $T(n) = T(n \operatorname{div} 2) + \Theta(n)$, следовательно, $T(n) = \Theta(n)$.

В качестве последней оптимизации можно избежать постоянного переувеличения длины, уточнив представление данных, а именно, заменив xs на пару $(\operatorname{length} xs, xs)$. Это приводит нас к окончательному варианту программы

$$\begin{array}{l} \operatorname{minfree} xs \\ \operatorname{minfrom} a (n, xs) \end{array} = \begin{array}{l} \operatorname{minfrom} 0 (\operatorname{length} xs, xs) \\ n == 0 \quad = a \\ m == b - a = \operatorname{minfrom} b (n - m, vs) \\ \text{otherwise} = \operatorname{minfrom} a (m, us) \\ \text{where } (us, vs) = \operatorname{partition} (<b) xs \\ \quad b = a + 1 + n \operatorname{div} 2 \\ \quad m = \operatorname{length} us \end{array}$$

Выясняется, что вышеприведённая программа примерно в два раза быстрее инкрементной программы на основе массива и на 20% быстрее, чем программа с использованием *accumArray*.

Заключительные замечания

Это была простая задача с как минимум двумя простыми решениями. Второе решение основывалось на известном методе проектирования алгоритмов, стратегии «разделяй и властвуй». Идея разбиения списка на те элементы, которые меньше заданного значения, и все остальные возникает во многих алгоритмах, например, в быстрой сортировке Хоара (Quicksort). При поиске алгоритма со сложностью $\Theta(n)$, использующего список из n

элементов, довольно заманчиво сразу обратиться к методу, обрабатывающему каждый элемент списка за константное или хотя бы амортизированное константное время. Однако рекурсивный процесс, который делает $\Theta(n)$ операций для сведения задачи к той же задаче не более чем половинного размера, также достаточно хорош.

Одно из различий между проектировщиками чисто функциональных и императивных алгоритмов состоит в том, что первые не предполагают существования массивов с операцией изменения за константное время, по крайней мере без некоторого объёма подготовительной работы. Для чисто функционального программиста операция изменения массива требует логарифмического по размеру массива времени¹. Это объясняет, почему иногда заметен логарифмический разрыв между функциональным и императивным решениями задачи. Но иногда, как здесь, этот разрыв при ближайшем рассмотрении исчезает.

¹По правде говоря, программисты-императивщики отлично знают, что константное время индексирования и изменения возможно только для маленьких массивов.

2

Превосходная задача

Введение

В этой жемчужине мы будем выполнять маленькое упражнение по программированию от Мартина Рема (Rem, 1988a). В то время как решение Рема использует бинарный поиск, наше будет ещё одним применением метода «разделяй и властвуй». Говорят, что некоторый элемент массива *превосходит* данный, если он больше и расположен правее, т.е. $x[j]$ превосходит $x[i]$, если $i < j$ и $x[i] < x[j]$. *Числом превосходства* (*surpasser count*) элемента массива называют количество элементов, его превосходящих. Например, вот числа превосходства для букв слова ТЕЛЕГРАФИСТ:

Т	Е	Л	Е	Г	Р	А	Ф	И	С	Т
1	6	4	5	5	3	4	0	2	1	0

Наибольшее число превосходства равно шести. Первое вхождение буквы Е имеет шесть превосходящих её элементов: буквы Л, Р, Ф, И, С и Т. Задача Рема заключается в вычислении наибольшего числа превосходства для массива длины $n > 1$ с помощью алгоритма со сложностью $O(n \log n)$.

Спецификация

Будем предполагать, что на входе вместо массива имеется список. Функция *msc* (сокращение от maximum surpasser count) может быть специфицирована следующим образом:

$$\begin{aligned}
msc &:: Ord\ a \Rightarrow [a] \rightarrow Int \\
msc\ xs &= maximum\ [scout\ z\ zs \mid z : zs \leftarrow tails\ xs] \\
scout\ x\ xs &= length\ (filter\ (x <) xs)
\end{aligned}$$

Значение $scout\ xs$ это число превосходства элемента x относительно списка xs , функция $tails$ возвращает непустые хвосты непустого списка в порядке убывания их длин¹:

$$\begin{aligned}
tails\ [] &= [] \\
tails\ (x : xs) &= (x : xs) : tails\ xs
\end{aligned}$$

Данное выше определение функции msc работает, но требует квадратичного времени.

Разделяй и властвуй

Учитывая заданную сложность $O(n \log n)$, кажется разумным обратиться к алгоритму по методу «разделяй и властвуй». Если нам удастся найти такую функцию $join$, что

$$msc\ (xs \uplus ys) = join\ (msc\ xs)\ (msc\ ys)$$

и которую можно вычислить за линейное время, то временная сложность $T(n)$ алгоритма «разделяй и властвуй» на списке длины n будет удовлетворять соотношению $T(n) = 2T(n/2) + O(n)$, а значит, $T(n) = O(n \log n)$. Однако довольно-таки очевидно, что такой функции $join$ не существует: единственное число $msc\ xs$ предоставляет слишком мало информации для любого подобного разбиения.

Наименьшим обобщением будет начать с таблицы всех чисел превосходства:

$$table\ xs = [(z, scout\ z\ zs) \mid z : zs \leftarrow tails\ xs]$$

Тогда $msc = maximum \cdot map\ snd \cdot table$. Теперь посмотрим, сможем ли мы найти линейную $join$, удовлетворяющую следующему:

$$table\ (xs \uplus ys) = join\ (table\ xs)\ (table\ ys)$$

Нам понадобится следующее свойство «разделяй и властвуй» для $tails$:

¹В отличие от стандартной функции языка Haskell с тем же именем, которая возвращает возможно пустые хвосты возможно пустого списка.

$$\mathit{tails} (xs \ ++ \ ys) = \mathit{map} (+ys) (\mathit{tails} \ xs) \ ++ \ \mathit{tails} \ ys$$

Проведём вычисления:

$$\begin{aligned} & \mathit{table} (xs \ ++ \ ys) \\ = & \quad \{ \text{определение} \} \\ & [(z, \mathit{scount} \ z \ zs) \mid z : zs \leftarrow \mathit{tails} (xs \ ++ \ ys)] \\ = & \quad \{ \text{свойство «разделяй и властвуй» для } \mathit{tails} \} \\ & [(z, \mathit{scount} \ z \ zs) \mid z : zs \leftarrow \mathit{map} (+ys) (\mathit{tails} \ xs) \ ++ \ \mathit{tails} \ ys] \\ = & \quad \{ \text{дистрибутивный закон для } \leftarrow \text{ по } ++ \} \\ & [(z, \mathit{scount} \ z \ (zs \ ++ \ ys)) \mid z : zs \leftarrow \mathit{tails} \ xs] \ ++ \\ & [(z, \mathit{scount} \ z \ zs) \mid z : zs \leftarrow \mathit{tails} \ ys] \\ = & \quad \{ \text{так как } \mathit{scount} \ z \ (zs \ ++ \ ys) = \mathit{scount} \ z \ zs + \mathit{scount} \ z \ ys \} \\ & [(z, \mathit{scount} \ z \ zs + \mathit{scount} \ z \ ys) \mid z : zs \leftarrow \mathit{tails} \ xs] \ ++ \\ & [(z, \mathit{scount} \ z \ zs) \mid z : zs \leftarrow \mathit{tails} \ ys] \\ = & \quad \{ \text{определение функции } \mathit{table} \text{ и } ys = \mathit{map} \ \mathit{fst} \ (\mathit{table} \ ys) \} \\ & [(z, c + \mathit{scount} \ z \ (\mathit{map} \ \mathit{fst} \ (\mathit{table} \ ys))) \mid (z, c) \leftarrow \mathit{table} \ xs] \ ++ \ \mathit{table} \ ys \end{aligned}$$

Следовательно, функцию *join* можно определить так:

$$\begin{aligned} \mathit{join} \ txs \ tys &= [(z, c + \mathit{tcount} \ z \ tys) \mid (z, c) \leftarrow txs] \ ++ \ tys \\ \mathit{tcount} \ z \ tys &= \mathit{scount} \ z \ (\mathit{map} \ \mathit{fst} \ tys) \end{aligned}$$

Однако проблема этого определения в том, что для вычисления *join txs tys* недостаточно линейного времени по длине *txs* и *tys*. Вычисление *tcount* можно было бы ускорить, если бы список *tys* был упорядочен по возрастанию первого компонента пары. В этом случае можно рассуждать так:

$$\begin{aligned} & \mathit{tcount} \ z \ tys \\ = & \quad \{ \text{определение } \mathit{tcount} \text{ и } \mathit{scount} \} \\ & \mathit{length} (\mathit{filter} (z <) (\mathit{map} \ \mathit{fst} \ tys)) \\ = & \quad \{ \text{так как } \mathit{filter} \ p \cdot \mathit{map} \ f = \mathit{map} \ f \cdot \mathit{filter} \ (p \cdot f) \} \\ & \mathit{length} (\mathit{map} \ \mathit{fst} \ (\mathit{filter} ((z <) \cdot \mathit{fst}) \ tys)) \\ = & \quad \{ \text{так как } \mathit{length} \cdot \mathit{map} \ f = \mathit{length} \} \\ & \mathit{length} (\mathit{filter} ((z <) \cdot \mathit{fst}) \ tys) \\ = & \quad \{ \text{так как } \mathit{tys} \text{ упорядочен по первому компоненту} \} \\ & \mathit{length} (\mathit{dropWhile} ((z \geq) \cdot \mathit{fst}) \ tys) \end{aligned}$$

Таким образом,

$$tcount\ z\ tys = length\ (dropWhile\ ((z \geq) \cdot fst)\ tys) \quad (2.1)$$

Это вычисление подсказывает, что было бы полезно поддерживать *table* в порядке возрастания первого компонента:

$$table\ xs = sort\ [(z, scount\ z\ zs) \mid z : zs \leftarrow tails\ xs]$$

Повторяя приведённое выше вычисление для упорядоченной версии *table*, получим, что

$$join\ txs\ tys = [(x, c + tcount\ x\ tys) \mid (x, c) \leftarrow txs] \mathbb{M}\ tys \quad (2.2)$$

где \mathbb{M} сливает два упорядоченных списка. Пользуясь этим определением, можно получить более эффективное рекурсивное определение *join*. Один из базовых случаев, $join\ []\ tys = tys$, очевиден. Другой, а именно, $join\ txs\ [] = txs$, следует из того, что $tcount\ x\ [] = 0$. Рекурсивную часть можно упростить, сравнивая x и y :

$$join\ txs@((x, c) : txs')\ tys@((y, d) : tys') \quad (2.3)$$

В языке Haskell символ $@$ вводит синоним, так что *txs* это синоним для $(x, c) : txs'$, аналогично для *tys*. Используя (2.2), выражение (2.3) можно свести к

$$((x, c + tcount\ x\ tys) : [(x, c + tcount\ x\ tys) \mid (x, c) \leftarrow txs']) \mathbb{M}\ tys$$

Чтобы узнать, какой элемент будет выдан операцией \mathbb{M} первым, нужно сравнить x и y . Если $x < y$, то это элемент слева и, поскольку согласно (2.1) $tcount\ x\ tys = length\ tys$, выражение (2.3) сводится к

$$(x, c + length\ tys) : join\ txs'\ tys$$

Если $x = y$, необходимо сравнить $c + tcount\ x\ tys$ и d . Но $d = tcount\ x\ tys'$ по определению *table*, а $tcount\ x\ tys = tcount\ x\ tys'$ согласно (2.1). Поэтому (2.3) сводится к $(y, d) : join\ txs\ tys'$. Такой же результат будет получен и в оставшемся случае $x > y$.

Собирая всё вместе и добавляя к *join* во избежание перевычисления длины дополнительный аргумент *length tys*, приходим к следующему алгоритму вычисления *table* на основе метода «разделяй и властвуй»:

$$\begin{aligned} table\ [x] &= [(x, 0)] \\ table\ xs &= join\ (m - n)\ (table\ ys)\ (table\ zs) \end{aligned}$$

where m = *length xs*
 n = $m \operatorname{div} 2$
 (ys, zs) = *splitAt n xs*

join 0 xs [] = *xs*

join n [] tys = *tys*

join n xs@((x, c) : xs') tys@((y, d) : tys')

| $x < y$ = $(x, c + n) : \operatorname{join} n \operatorname{xs}' \operatorname{tys}'$

| $x \geq y$ = $(y, d) : \operatorname{join} (n - 1) \operatorname{xs} \operatorname{tys}'$

Так как *join* выполняется за линейное время, список *table* вычисляется за $O(n \log n)$ операций, а значит такова же сложность *msc*.

Заключительные замечания

Невозможно вычислить *table* с помощью алгоритма, более быстрого, чем $O(\log n)$. Причина в том, что если *xs* это список без повторяющихся элементов, то *table xs* предоставляет достаточно информации для определения такой сортирующей перестановки списка *xs*. Более того, никакие последующие сравнения между элементами списка не требуются. Фактически *table xs* соотносится с таблицей инверсий для перестановки из n элементов; см. Кнута (Knuth, 1998): *table xs* это просто таблица инверсий для *reverse xs*. Так как основанная на сравнениях сортировка n элементов требует $\Omega(n \log n)$ операций, то столько же требует и вычисление *table*.

Как было сказано во введении, решение Рема (Rem, 1988b) отличается тем, что оно использует итеративный алгоритм с применением бинарного поиска. Программист-императивщик также мог бы обратиться к алгоритму «разделяй и властвуй», но он, вероятно, предпочёл бы алгоритм с обработкой массива просто потому, что он требует меньше памяти.

Литература

- Knuth, D. E. (1998). *The Art of Computer Programming, Volume 3: Sorting and Searching*, second edition. Reading, MA: Addison-Wesley. [Имеется русский перевод: Кнут Д. Искусство программирования, том 3. Сортировка и поиск. 2-е изд. М.: Вильямс. 2007.]
- Rem, M. (1988a). Small programming exercises 20. *Science of Computer Programming* **10** (1), 99–105.
- Rem, M. (1988b). Small programming exercises 21. *Science of Computer Programming* **10** (3), 319–25.

3

Улучшаем седловой поиск

Действие происходит на занятии по проектированию функциональных алгоритмов. В классе четверо студентов: Анна, Иван, Мария и Фёдор.

Учитель. Доброе утро, ребята. Сегодня я бы хотел, чтобы вы спроектировали функцию *invert* с двумя параметрами: функция *f*, действующая из множества пар натуральных чисел во множество натуральных чисел, и натуральное число *z*. Значение *invert f z* должно быть списком всех пар (x, y) , удовлетворяющих равенству $f(x, y) = z$. Вы можете предполагать только то, что *f* строго возрастает по каждому из аргументов, и ничего иного.

Иван. Кажется, это простая задача. Так как *f* действует на натуральных числах и возрастает по каждому из аргументов, то мы знаем, что из равенства $f(x, y) = x$ следует, что $x \leq z$ и $y \leq z$. Поэтому можно определить *invert* простым поиском всех возможных пар значений:

$$\text{invert } f \ z = [(x, y) \mid x \leftarrow [0..z], y \leftarrow [0..z], f(x, y) == z]$$

Разве это не решение?

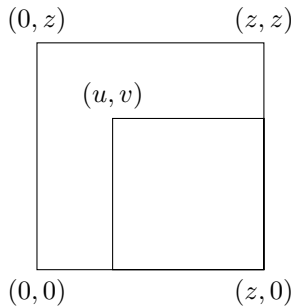
Учитель. Да, это решение, но оно требует $(z + 1)^2$ вычислений значения функции *f*. Поскольку вычисление значения *f* может оказаться очень долгим, я бы хотел увидеть решение, требующее как можно меньшее количество вызовов *f*.

Фёдор. Ну, нетрудно уменьшить количество вычислений вдвое. Так как при условии, что f возрастает, $f(x, y) \geq x + y$, то поиск можно проводить только по тем значениям, которые расположены на или ниже диагонали квадрата:

$$\text{invert } f z = [(x, y) \mid x \leftarrow [0..z], y \leftarrow [0..z - x], f(x, y) == z]$$

А если задуматься, то обе верхние границы можно вообще заменить на $z - f(0, 0)$ и $z - x - f(0, 0)$. Тогда, если $z < f(0, 0)$, то поиск сразу останавливается.

Анна. В предположении, что порядок найденных решений не важен, я думаю, результат можно ещё улучшить. Метод Ивана выполняет поиск от начала координат в левом нижнем углу квадрата размера $z + 1$, обходя его столбец за столбцом. Мы же можем ускориться, если начнём в левом верхнем углу квадрата в точке $(0, z)$. На каждом этапе область поиска ограничена прямоугольником с левым верхним углом в точке (u, v) и правым нижним углом в точке $(z, 0)$. Вот картинка:



Давайте определим функцию

$$\text{find } (u, v) f z = [(x, y) \mid x \leftarrow [u..z], y \leftarrow [v, v - 1..0], f(x, y) == z]$$

Таким образом, $\text{invert } f z = \text{find } (0, z) f z$. Теперь нетрудно найти более эффективную реализацию функции find .

Во-первых, если $u > z$ или $v < 0$, то очевидно, $\text{find } (u, v) f z = []$. В противном случае будем рассматривать различные возможные значения $f(u, v)$. Если $f(u, v) < z$, то остаток столбца u можно отбросить, так как при $v' < v$ имеет место неравенство $f(u, v') < f(u, v) < z$. Если $f(u, v) > z$, то аналогично отбрасываем остаток строки v . Наконец, если $f(u, v) = z$, то можно запомнить (u, v) и отбросить как столбец u , так и строку v .

Вот улучшенная версия функции *invert*:

```
invert f z = find (0, z) f z
find (u, v) f z
  | u > z ∨ v < 0 = []
  | z' < z       = find (u + 1, v) f z
  | z' == z      = (u, v) : find (u + 1, v - 1) f z
  | z' > z       = find (u, v - 1) f z
  where z' = f(u, v)
```

В худшем случае, если *find* обходит периметр квадрата от левого верхнего угла до правого нижнего, она делает $2z + 1$ вызовов функции *f*. В лучшем случае, если *find* идёт напрямую либо к нижней, либо к правой границе, ей требуется лишь $z + 1$ вызов.

Фёдор. Ты можешь сократить область поиска ещё сильнее, потому что исходный квадрат с левым верхним углом $(0, z)$ и правым нижним углом $(z, 0)$ является слишком избыточной оценкой области нахождения требуемого значения. Предположим, что мы предварительно вычислили такие m и n , что

```
m = maximum (filter (λy → f(0, y) ≤ z) [0..z])
n = maximum (filter (λx → f(x, 0) ≤ z) [0..z])
```

Тогда можно определить $invert\ f\ z = find(0, m)\ f\ z$, где *find* имеет в точности ту же форму, в какой её определила Анна, за тем исключением, что первое охранное выражение становится таким: $u > n \vee v < 0$. Другими словами, вместо того, чтобы анализировать весь квадрат $(z + 1) \times (z + 1)$, можно ограничиться квадратом размера $(m + 1) \times (n + 1)$.

Решающее соображение состоит в том, что m и n можно вычислить бинарным поиском. Пусть g — возрастающая функция на множестве натуральных чисел, предположим также, что x , y и z удовлетворяют условию $g\ z \leq z < g\ y$. Для отыскания единственного значения m , где $m = bsearch\ g\ (x, y)\ z$, на отрезке $x \leq m < y$ и с учётом $g\ m \leq z < g\ (m + 1)$ можно использовать инварианты $g\ a \leq z < g\ b$ и $x \leq a < b \leq y$. Таким образом, получаем функцию

```
bsearch g (a, b) z
  | a + 1 == b = a
  | g m ≤ z   = bsearch g (m, b) z
  | otherwise = bsearch g (a, m) z
  where m = (a + b) div 2
```

Так как $a + 1 < b \Rightarrow a < m < y$, то ни $g x$, ни $g y$ алгоритмом не вычисляются, поэтому они могут быть фиктивными значениями. В частности, имеем

$$\begin{aligned} m &= \text{bsearch}(\lambda y \rightarrow f(0, y))(-1, z + 1) z \\ n &= \text{bsearch}(\lambda x \rightarrow f(x, 0))(-1, z + 1) z \end{aligned}$$

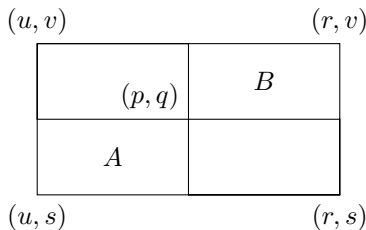
где f дополнена фиктивными значениями $f(0, -1) = 0$ и $f(-1, 0) = 0$.

Этой версии *invert* требуется примерно $2 \log z + m + n$ вызовов f в худшем случае и $2 \log z + m \min n$ в лучшем. Так как m и n могут оказаться значительно меньше z , к примеру, если $f(x, y) = 2^x + 3^y$, то мы получаем алгоритм, требующий в худшем случае только $O(\log z)$ операций.

Учитель. Мои поздравления, Анна и Фёдор, вы переоткрыли важную стратегию поиска, названную Дэвидом Грисом (David Gries) *седловым поиском*; см. работы Бэкхауза (Backhouse, 1986), Дейкстры (Dijkstra, 1985) и Гриса (Gries, 1981). Мне представляется, что Грис назвал его так, потому что форма трёхмерного графика f с наименьшим элементом в левом нижнем углу, наибольшим в правом верхнем и двумя крыльями напоминает седло. Основная идея, как заметила Анна, в том, чтобы начать поиск на кончике одного из крыльев, а не в точках с наименьшим или наибольшим значением. Рассматривая эту задачу, Дейкстра (Dijkstra, 1985) также указал на преимущества использования логарифмического поиска для определения подходящего начального прямоугольника.

Мария. А что если мы попробуем решение по методу «разделяй и властвуй»? В смысле, почему бы в начале не смотреть на элемент в середине прямоугольника? Уверена, что имеет смысл применить двумерный аналог бинарного поиска.

Предположим, мы ограничили поиск прямоугольником с левым верхним углом (u, v) и правым нижним (r, s) . Почему бы вместо того, чтобы смотреть на $f(u, v)$, не проанализировать $f(p, q)$, где $p = (u + r) \operatorname{div} 2$, а $q = (v + s) \operatorname{div} 2$? Вот картинка:



Если $f(p, q) < z$, то мы можем отбросить все элементы левого нижнего прямоугольника A . Аналогично, если $f(p, q) > z$, можно отбросить правый верхний прямоугольник B . Если же $f(p, q) = z$, то можно отбросить оба.

Я знаю, что эта стратегия нарушает свойство, предложенное Анной, о том, что область поиска всегда прямоугольник. Вместо этого получаем либо два прямоугольника, либо Γ -образную форму. Но мы же функциональные программисты и можем не ограничивать себя простыми циклами: алгоритм «разделяй и властвуй» столь же лёгок в реализации, что и итеративный, потому что оба можно выразить рекурсией.

Иван. Тебе всё равно придётся иметь дело с Γ -образной формой. Конечно, её можно разбить на два прямоугольника. Фактически это можно сделать двумя способами: либо горизонтальным разрезом, либо вертикальным. Позволь мне дать грубую оценку. Рассмотрим прямоугольник размера $m \times n$ и обозначим $T(m, n)$ количество требуемых для поиска по нему вычислений f . Если $m = 0$ или $n = 0$, то искать нечего. Если $m = 1$ или $n = 1$, имеем

$$\begin{aligned} T(1, n) &= 1 + T(1, \lceil n/2 \rceil) \\ T(m, 1) &= 1 + T(\lceil m/2 \rceil, 1) \end{aligned}$$

В противном случае, если $m \geq 2$ и $n \geq 2$, можно отбросить прямоугольник размера, как минимум, $\lfloor m/2 \rfloor \times \lfloor n/2 \rfloor$. Если мы делаем горизонтальный разрез, то остаётся два прямоугольника: один размера $\lfloor m/2 \rfloor \times \lceil n/2 \rceil$, а второй $\lceil m/2 \rceil \times n$. Таким образом,

$$T(m, n) = 1 + T(\lfloor m/2 \rfloor, \lceil n/2 \rceil) + T(\lceil m/2 \rceil, n)$$

Если же мы делаем вертикальный разрез, то имеем

$$T(m, n) = 1 + T(\lceil m/2 \rceil, \lfloor n/2 \rfloor) + T(m, \lfloor n/2 \rfloor)$$

Я, правда, не вижу решения этих рекуррентных соотношений.

Фёдор. Если сделать одновременно и горизонтальный, и вертикальный разрезы, то останется три прямоугольника, так что при $m \geq 2$ и $n \geq 2$ будем иметь

$$T(m, n) = 1 + T(\lceil m/2 \rceil, \lfloor n/2 \rfloor) + T(\lceil m/2 \rceil, \lceil n/2 \rceil) + T(\lfloor m/2 \rfloor, \lceil n/2 \rceil)$$

Я могу решить это соотношение. Положим $U(i, j) = T(2^i, 2^j)$, теперь

$$\begin{aligned} U(i, 0) &= i \\ U(0, j) &= j \\ U(i+1, j+1) &= 1 + 3U(i, j) \end{aligned}$$

Решением будет $U(i, j) = 3^k(|j - i| + 1/2) - 1/2$, где $k = i \min j$, это можно проверить индукцией. Следовательно, если $m \leq n$, то

$$T(m, n) \leq 3^{\log m} \log(2n/m) = m^{1.59} \log(2n/m)$$

Это лучше, чем $m + n$, при условии, что m значительно меньше n .

Иван. Я не думаю, что решение с тремя прямоугольниками столь же хорошее, что и с двумя. Следуя твоему подходу, Фёдор, давай положим $U(i, j) = T(2^i, 2^j)$. В предположении, что $i \leq j$, выполняя горизонтальный разрез, получаем

$$\begin{aligned} U(0, j) &= j \\ U(i+1, j+1) &= 1 + U(i, j) + U(i, j+1) \end{aligned}$$

Решением является $U(i, j) = 2^i(j - i/2 + 1) - 1$, что также можно проверить индукцией. Поэтому

$$T(m, n) \leq m \log(2n/\sqrt{m})$$

Если $i \geq j$, то следует делать вертикальный разрез вместо горизонтального. Тогда мы получим алгоритм, требующий самое большее $n \log(2m/\sqrt{n})$ вызовов f . В любом случае, если один из параметров m или n значительно меньше другого, то мы имеем алгоритм лучше, чем при седловом поиске.

Анна. Пока вы двое решали рекуррентные соотношения, я думала над проблемой нижней границы сложности для функции *invert*. Рассмотрим различные возможные результаты при поиске в прямоугольнике размера $m \times n$. Предположим, что имеется $A(m, n)$ различных решений. Каждая проверка $f(x, y)$ относительно z допускает три возможных результата, поэтому высота h тернарного дерева проверок должна удовлетворять условию $h \geq \log_3 A(m, n)$. Считая, что $A(m, n)$ достижимо, получаем нижнюю границу для числа проверок, которые необходимо выполнить. Ситуация похожа на сортировку n элементов двоичными сравнениями: имеется $n!$ возможных результатов, поэтому любой алгоритм сортировки должен в худшем случае провести как минимум $\log_2 n!$ сравнений.

$A(m, n)$ достичь нетрудно: каждый список пар (x, y) , где $0 \leq x < n$ и $0 \leq y < m$, удовлетворяющих $f(x, y) = z$, находится в соотношении один к одному со ступенчатой функцией, график которой идёт из левого верхнего угла прямоугольника размера $m \times n$ в его правый нижний угол. Значение z появляется на внутренних углах ступенек. Разумеется, этот график

необязательно задаёт именно тот путь, по которому идёт функция *find*. Количество таких путей $\binom{m+n}{n}$, а значит, таким же будет значение $A(m, n)$.

По другому этот результат можно получить, если предположить, что имеется k решений. Значение z можно разместить в m строках k раз ровно $\binom{m}{k}$ способами, и всякий раз будет $\binom{n}{k}$ возможных вариантов для столбцов. Следовательно,

$$A(m, n) = \sum_{k=0}^m \binom{m}{k} \binom{n}{k} = \binom{m+n}{n}$$

поскольку это суммирование является примером свёртки Вандермонда; см. (Graham, 1989b). Беря логарифмы, получаем нижнюю границу

$$\log A(m, n) = \Omega(m \log(1 + n/m) + n \log(1 + m/n))$$

Это приближение показывает, что если $m = n$, то не удастся сделать менее, чем $\Omega(m+n)$ шагов. Но если $m \leq n$, то $m \leq n \log(1 + m/n)$, поскольку $x \leq \log(1+x)$ при $0 \leq x \leq 1$. Таким образом, $A(m, n) = \Omega(m \log(n/m))$. Решение Ивана не слишком близко к этой границе, потому что в его алгоритме в случае $m \leq n$ число операций оценивается как $O(m \log(n/\sqrt{m}))$.

Мария. Я не думаю, что в решении Ивана действительно необходимо применять метод «разделяй и властвуй». Есть и другие способы использовать бинарный поиск. Например, можно просто выполнить m бинарных поисков, по одному для каждой строки. Так мы получим решение со сложностью $O(m \log n)$. Но мне кажется, что её можно улучшить, достигнув оптимальную асимптотическую сложность $O(m \log(n/m))$ в предположении, что $m \leq n$.

Допустим, как и прежде, что поиск ограничен прямоугольником с левым верхним углом (u, v) и правым нижним (r, s) . Таким образом, имеется $r - u$ столбцов и $s - v$ строк. Далее, предположим, что $v - s \leq r - u$, т.е. столбцов как минимум столько же, сколько и строк. Пусть мы выполняем бинарный поиск по средней строке, $q = (v + s) \text{ div } 2$, и ищем такое p , что $f(p, q) \leq z < f(p + 1, q)$. Если $f(p, q) < z$, то нужно продолжать поиск только в двух прямоугольниках $((u, v), (p, q + 1))$ и $((p + 1, q - 1), (r, s))$. Если $f(p, q) = z$, то можно вырезать столбец p и продолжить поиск в прямоугольниках $((u, v), (p - 1, q + 1))$ и $((p + 1, q - 1), (r, s))$. Рассуждения в случае, когда строк больше, чем столбцов, аналогичны. В результате, можно отбросить половину элементов массива с помощью логарифмического числа испытаний.

Вот алгоритм, который я имею в виду: реализуем *invert*

```

find (u, v) (r, s) f z
  | u > r ∨ v < s = []
  | v - s ≤ r - u = rfind (bsearch (λx → f(x, q)) (u - 1, r + 1) z)
  | otherwise     = cfind (bsearch (λy → f(p, y)) (s - 1, v + 1) z)
where
  p      = (u + r) div 2
  q      = (v + s) div 2
  rfind p = (if f(p, q) == z then (p, q) : find (u, v) (p - 1, q + 1) f z
            else find (u, v) (p, q + 1) f z) ++
            find (p + 1, q - 1) (r, s) f z
  cfind q = find (u, v) (p - 1, q + 1) f z ++
            (if f(p, q) == z then (p, q) : find (p + 1, q - 1) (r, s) f z
            else find (p + 1, q) (r, s) f z)

```

Рис. 3.1: пересмотренное определение функции *find*

```

invert f z = find (0, m) (n, 0) f z
  where m = bsearch (λy → f(0, y)) (-1, z + 1) z
        n = bsearch (λx → f(x, 0)) (-1, z + 1) z

```

где $find (u, v) (r, s) f z$, приведённая на рис. 3.1, выполняет поиск по прямоугольнику с левым верхним углом (u, v) и правым нижним углом (r, s) .

Что же касается анализа, пусть $T(m, n)$ снова обозначает количество вызовов, необходимых для поиска по прямоугольнику $m \times n$. Предположим, что $m \leq n$. В лучшем случае, когда каждый бинарный поиск по строке возвращает самый левый или самый правый элемент, имеем $T(m, n) = \log n + T(m/2, n)$, т.е. $T(m, n) = O(\log m \times \log n)$. В худшем случае, когда бинарный поиск возвращает средний элемент, получается

$$T(m, n) = \log n + 2T(m/2, n/2)$$

Чтобы решить это соотношение, положим $U(i, j) = T(2^i, 2^j)$. Тогда имеем

$$U(i, j) = \sum_{k=0}^{i-1} 2^k (j - k) = O(2^i (j - i))$$

Следовательно, $T(m, n) = O(m \log(n/m))$, что асимптотически оптимально согласно полученной Анной нижней оценке.

Учитель. Отличная работа, все четверо! Удивительно, что примерно за 25 лет, в которые седловой поиск использовался в качестве примера формального построения программ, никто, как кажется, не заметил, что это алгоритм с не самой лучшей асимптотикой.

Послесловие

Реальная история, стоящая за этой жемчужиной, такова. Я решил использовать седловой поиск в качестве упражнения для кандидатов, проходящих собеседование на поступление в Оксфорд. Им давался двумерный числовой массив, возрастающий вдоль каждой строки и каждого столбца, и предлагалось найти систематический способ обнаружения всех вхождений заданного числа. Я хотел убедить их, что поиск из левого верхнего угла в правый нижний является хорошей стратегией. Однако те кандидаты, который изучали в школе информатику, пытались использовать бинарный поиск либо из середины каждой строки, либо из центра прямоугольника. Полагая, что седловой поиск является золотым стандартом для решения данной задачи, я намекал им, что следует поменять ход мысли. Только спустя некоторое время я заинтересовался тем, могли ли они быть правы.

Помимо описания нового алгоритма для старой задачи, как мне кажется, есть два других методологических аспекта, достойных упоминания. Во-первых, формальное построение программ испытывает сильное влияние со стороны доступных в целевом языке программирования методов вычисления. Последнее решение, приведённое Марией, вряд ли кто-то назовёт элегантным, тем не менее, оно достаточно элементарно в ситуации, когда рекурсия и конкатенация списков являются базовыми операциями, однако его реализация в языке с массивами и циклами может оказаться довольно затруднительной. Во-вторых, как совершенно уверены проектировщики алгоритмов, формальное построение программ должно сопровождаться идеями относительно возможных путей улучшения эффективности. Такие идеи могут возникнуть частично благодаря решению рекуррентных соотношений, а частично из определения нижних границ.

Впервые эта жемчужина появилась в (Bird, 2006). Один из рецензентов исходной работы указал следующее:

Сложная реализация помимо собственно эффективности приносит и некоторые накладные расходы, которые так часто упускают из виду при анализе, подобном представленному в этой работе. Если автор действительно хочет убедить нас, что