

1

The smallest free number

Introduction

Consider the problem of computing the smallest natural number not in a given finite set X of natural numbers. The problem is a simplification of a common programming task in which the numbers name objects and X is the set of objects currently in use. The task is to find some object not in use, say the one with the smallest name.

The solution to the problem depends, of course, on how X is represented. If X is given as a list without duplicates and in increasing order, then the solution is straightforward: simply look for the first gap. But suppose X is given as a list of distinct numbers in no particular order. For example,

[08, 23, 09, 00, 12, 11, 01, 10, 13, 07, 41, 04, 14, 21, 05, 17, 03, 19, 02, 06]

How would you find the smallest number not in this list?

It is not immediately clear that there is a linear-time solution to the problem; after all, sorting an arbitrary list of numbers cannot be done in linear time. Nevertheless, linear-time solutions do exist and the aim of this pearl is to describe two of them: one is based on Haskell arrays and the other on divide and conquer.

An array-based solution

The problem can be specified as a function *minfree*, defined by

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

The expression $us \setminus vs$ denotes the list of those elements of us that remain after removing any elements in vs :

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

The function *minfree* is executable but requires $\Theta(n^2)$ steps on a list of length n in the worst case. For example, evaluating *minfree* $[n-1, n-2 \dots 0]$ requires evaluating $i \notin [n-1, n-2 \dots 0]$ for $0 \leq i \leq n$, and thus $n(n+1)/2$ equality tests.

The key fact for both the array-based and divide and conquer solutions is that not every number in the range $[0 \dots \text{length } xs]$ can be in xs . Thus the smallest number not in xs is the smallest number not in *filter* $(\leq n) xs$, where $n = \text{length } xs$. The array-based program uses this fact to build a checklist of those numbers present in *filter* $(\leq n) xs$. The checklist is a Boolean array with $n+1$ slots, numbered from 0 to n , whose initial entries are everywhere *False*. For each element x in xs and provided $x \leq n$ we set the array element at position x to *True*. The smallest free number is then found as the position of the first *False* entry. Thus, *minfree* = *search* · *checklist*, where

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

The function *search* takes an array of Booleans, converts the array into a list of Booleans and returns the length of the longest initial segment consisting of *True* entries. This number will be the position of the first *False* entry.

One way to implement *checklist* in linear time is to use the function *accumArray* in the Haskell library *Data.Array*. This function has the rather daunting type

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

The type constraint $Ix \ i$ restricts i to be an *Index* type, such as *Int* or *Char*, for naming the indices or positions of the array. The first argument is an “accumulating” function for transforming array entries and values into new entries, the second argument is an initial entry for each index, the third argument is a pair naming the lower and upper indices and the fourth is an association list of index–value pairs. The function *accumArray* builds an array by processing the association list from left to right, combining entries and values into new entries using the accumulating function. This process takes linear time in the length of the association list, assuming the accumulating function takes constant time.

The function *checklist* is defined as an instance of *accumArray*:

$$\begin{aligned} \text{checklist} &:: [\text{Int}] \rightarrow \text{Array Int Bool} \\ \text{checklist } xs &= \text{accumArray } (\vee) \text{ False } (0, n) \\ &\quad (\text{zip } (\text{filter } (\leq n) xs) (\text{repeat True})) \\ &\quad \textbf{where } n = \text{length } xs \end{aligned}$$

This implementation does not require the elements of xs to be distinct, but does require them to be natural numbers.

It is worth mentioning that *accumArray* can be used to sort a list of numbers in linear time, provided the elements of the list all lie in some known range $(0, n)$. We replace *checklist* by *countlist*, where

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

Then $\text{sort } xs = \text{concat} [\text{replicate } k \ x \mid (x, k) \leftarrow \text{countlist } xs]$. In fact, if we use *countlist* instead of *checklist*, then we can implement *minfree* as the position of the first 0 entry.

The above implementation builds the array in one go using a clever library function. A more prosaic way to implement *checklist* is to tick off entries step by step using a constant-time update operation. This is possible in Haskell only if the necessary array operations are executed in a suitable monad, such as the state monad. The following program for *checklist* makes use of the library *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
```

This solution would not satisfy the pure functional programmer because it is essentially a procedural program in functional clothing.

A divide and conquer solution

Now we turn to a divide and conquer solution to the problem. The idea is to express $\text{minfree } (xs \mathbin{++} ys)$ in terms of $\text{minfree } xs$ and $\text{minfree } ys$. We begin by recording the following properties of \setminus :

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

These properties reflect similar laws about sets in which set union \cup replaces $\mathbin{++}$ and set difference \setminus replaces \setminus . Suppose now that as and vs are disjoint, meaning $as \setminus vs = as$, and that bs and us are also disjoint, so $bs \setminus us = bs$. It follows from these properties of $\mathbin{++}$ and \setminus that

$$(as \mathbin{++} bs) \setminus (us \mathbin{++} vs) = (as \setminus us) \mathbin{++} (bs \setminus vs)$$

Now, choose any natural number b and let $as = [0 .. b-1]$ and $bs = [b..]$. Furthermore, let $us = \text{filter } (< b) \text{ } xs$ and $vs = \text{filter } (\geq b) \text{ } xs$. Then as and vs are disjoint, and so are bs and us . Hence

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

Haskell provides an efficient implementation of a function *partition* p that partitions a list into those elements that satisfy p and those that do not. Since

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

we obtain, still for any natural number b , that

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

The next question is: can we implement the test $\text{null } ([0 .. b-1] \setminus us)$ more efficiently than by direct computation, which takes quadratic time in the length of us ? Yes, the input is a list of distinct natural numbers, so is us . And every element of us is less than b . Hence

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

Note that the array-based solution did not depend on the assumption that the given list did not contain duplicates, but it is a crucial one for an efficient divide and conquer solution.

Further inspection of the above code for *minfree* suggests that we should generalise *minfree* to a function, *minfrom* say, defined by

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

where every element of xs is assumed to be greater than or equal to a . Then, provided b is chosen so that both $\text{length } us$ and $\text{length } vs$ are less than $\text{length } xs$, the following recursive definition of *minfree* is well-founded:

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

It remains to choose b . Clearly, we want $b > a$. And we would also like to choose b so that the maximum of the lengths of us and vs is as small as possible. The right choice of b to satisfy these requirements is

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

where $n = \text{length } xs$. If $n \neq 0$ and $\text{length } us < b - a$, then

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

And, if $\text{length } us = b - a$, then

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

With this choice the number of steps $T(n)$ for evaluating $\text{minfrom } 0 \ xs$ when $n = \text{length } xs$ satisfies $T(n) = T(n \operatorname{div} 2) + \Theta(n)$, with the solution $T(n) = \Theta(n)$.

As a final optimisation we can avoid repeatedly computing length with a simple data refinement, representing xs by a pair $(\text{length } xs, xs)$. That leads to the final program

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

It turns out that the above program is about twice as fast as the incremental array-based program, and about 20% faster than the one using *accumArray*.

Final remarks

This was a simple problem with at least two simple solutions. The second solution was based on a common technique of algorithm design, namely divide and conquer. The idea of partitioning a list into those elements less than a given value, and the rest, arises in a number of algorithms, most notably Quicksort. When seeking a $\Theta(n)$ algorithm involving a list of n elements, it is tempting to head at once for a method that processes each element of the list in constant time, or at least in amortized constant time. But a recursive process that performs $\Theta(n)$ processing steps in order to reduce the problem to another instance of at most half the size is also good enough.

One of the differences between a pure functional algorithm designer and a procedural one is that the former does not assume the existence of arrays with a constant-time update operation, at least not without a certain amount of plumbing. For a pure functional programmer, an update operation takes logarithmic time in the size of the array.¹ That explains why there sometimes seems to be a logarithmic gap between the best functional and procedural solutions to a problem. But sometimes, as here, the gap vanishes on a closer inspection.

¹ To be fair, procedural programmers also appreciate that constant-time indexing and updating are only possible when the arrays are small.