

Section 2.5 (2.6, 6th edition): Analysis of Algorithms

September 19, 2007

Abstract

By **analysis of algorithms** we mean the study of the efficiency of the algorithms. In this section we will measure the efficiency of an algorithm by counting operations (and of course we are generally shooting for a **small** number, in our endless pursuit of optimization).

1 Counting operations directly

In algorithm *SequentialSearch* (p. 148/170), we search for element x in a list of n items. *SequentialSearch* is a direct method, by comparison with algorithm *BinarySearch* (p. 130/138), which is recursive. Is one algorithm more efficient than the other?

In the *SequentialSearch*, there are three rather interesting cases:

- we find x on the very first try (total comparisons: $1!$). This is called the “best-case” scenario.
- we find x on the last try (total comparisons: n). This is the “worst-case” scenario.
- On average, we require $(n+1)/2$ comparisons, remembering Gauss: we sum up all the cases from 1 to n , and divide by n :

$$\frac{1}{n} \sum_{i=1}^n i = \frac{1}{n} \frac{n(n+1)}{2} = \frac{n+1}{2}$$

We will consider the worst-case scenario as the benchmark.

2 Counting Using Recurrence Relations

Algorithm *BinarySearch* is recursive: it calls itself. Starting from a list of length n it makes one comparison and then calls itself with a list of half its initial length. Hence the number of comparisons for the list of length n , $C(n)$, would be (in the worst case)

$$C(n) = C(\text{floor}(n/2)) + 1$$

and $C(1) = 1$. That floor function is a pain, but is necessary since n may be odd.

Forgetting the floor for the moment, use the “expand, guess, and verify” approach: in the worst-case scenario, the algorithm will find the element (or not) on its last check (when it’s down to a list of length 1).

$$C(n) = C(n/2) + 1 = (C(n/4) + 1) + 1 = ((C(n/8) + 1) + 1) + 1 = \dots$$

Obviously this is only going to work easily (in the sense that $C(n/8)$, etc., make sense) if n is a power of 2. Assume therefore that $n = 2^m$, for integer m . This allows us to throw away the floor function, and makes all quotients reasonable.

Consider a change of variable: in

$$C(2^m) = C(2^{m-1}) + 1$$

we define $T(m) = C(2^m)$ (think of T as a composition of functions, C and 2^x); hence

$$T(m) = T(m-1) + 1$$

Note that $T(0) = C(1) = 1$. We can solve easily to get a closed-form solution of

$$T(m) = m + 1$$

Hence, $C(n) = C(2^m) = T(m) = m + 1 = \log_2(n) + 1$. This compares quite favorably with the worst-case estimate from *SequentialSearch*, which would be n (linear in n).

$$n = 2^m \\ \log_2 n = m \log_2 2$$

(For those of you who've forgotten, the log function grows much more slowly than a linear function.)

Let's look at the general recurrence relation of the "divide and conquer" variety: given

$$\begin{aligned} S(1) &= a \\ S(n) &= cS(n/2) + g(n) \end{aligned}$$

Assume $n = 2^m$ for some integer m . Then

$$\begin{aligned} S(2^0) &= a \\ S(2^m) &= cS(2^{m-1}) + g(2^m) \end{aligned}$$

Now we perform the change of variables: let $T(m) = S(2^m)$, so that

$$\begin{aligned} T(0) &= a \\ T(m) &= cT(m-1) + g(2^m) \end{aligned}$$

Using formula (8) of section 2.4/2.5, p. 134/150, we get

$$T(m) = c^{m-1}T(1) + \sum_{i=2}^m c^{m-i}g(2^i)$$

Then reindexing, since we start with 0 rather than 1, we get

$$T(m) = c^mT(0) + \sum_{i=1}^m c^{m-i}g(2^i)$$

Finally, substituting back in S and n , we get

$$S(2^m) = c^{\log_2 n} a + \sum_{i=1}^{\log_2 n} c^{\log_2 n - i} g(2^i)$$

$$\begin{aligned} S(n) &= c S(n/2) \\ &+ g(n) \end{aligned}$$

Whew!

The *BinarySearch* algorithm starts with a sorted list, which is not a requirement for the *SequentialSearch* algorithm; so the comparison isn't really fair. What if we add a sort?

Example: Exercise 13/18, p. 156/179

Example: Exercise 14/19, p. 156/179

Example: Exercise 15/20, p. 156/179

Write a recurrence relation $n = 2^m$

Sorting a list of length 1: $C(1) = 0$

$$C(n) = 2 C(n/2) + \underline{n-1}$$

merging two halves of length $n/2$

Example: Exercise 16/21, p. 156/179

$$c = 2$$

$$g(n) = n - 1$$

$$a = 0$$

$$\begin{aligned} C(n) &= c \cancel{\log n} g(1) + \sum_{i=1}^{\log n} c^{\log n - i} g(2^i) \\ &= \sum_{i=1}^{\log n} 2^{\log n - i} (2^i - 1) = \end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^{\log_2 n} 2^{\log_2 n - i} - (2^{\log_2 n - i}) \\
&= \log_2 n \cdot 2^{\log_2 n} - \sum_{i=1}^{\log_2 n} 2^{\log_2 n - i}
\end{aligned}$$

So we can carry out the *BinarySearch* algorithm following a *MergeSort* (see the exercises above for its definition), with

$$\log_2(n) + 1 + n \log_2(n) - n + 1$$

or

$$(n + 1) \log_2(n) + 2$$

operations, compared with n operations for *SequentialSearch* - which wins in this case! $(n + 1) \log_2(n)$ is *superlinear* - grows faster than the linear function n .

If we had started with a sorted list, however, it would make no sense to use *SequentialSearch*, since *BinarySearch* is so much more efficient in that case.

Also, if we are doing multiple searches with the same lists, then the costs of not sorting begin to add up. The sorting is an initial (or fixed) cost; then there is a benefit each time one sorts thereafter for the merge-sort. Hence we can compute the cut-off value of the number of searches for a given list size which would justify sorting and then using the binary search algorithm.

3 Other criteria

An algorithm should not be analyzed quite so one-dimensionally as we've done here, of course: there may be other issues (such as how easily parallelized an algorithm is, for example) which are more important than simple operation counts.

As demonstrated in the case of the Euclidean Algorithm (or gcd) in this section, we may simply be shooting for an upper bound on the number of operations required (even worse than the worst case scenario!), when actual worst-case numbers are hard to come by.

The Euclidean Algorithm finds the greatest common denominator of two integers (See p. 114/122).

Actually, in this case, worst-case numbers are easy to get: the worst case for the Euclidean algorithm is a pair of consecutive Fibonacci numbers (there they are again, those rascals!).

This is investigated in problems 20-23/25-28. An example pair of consecutive Fibonacci numbers would be 5 and 3, or 89 and 55.

$$\begin{aligned}
&= \log_2 n \cdot 2^{\log_2 n} - \sum_{i=1}^{\log_2 n} 2^{\log_2 n - i} \\
&= \underbrace{\hspace{10em}}_{11} - \sum_{i=0}^{\log_2 n - 1} 2^i
\end{aligned}$$

=

$$\frac{2^0 + 2^1 + \dots + 2^{\log_2 n - 1}}{1 + r + r^2 + \dots + r^{\log_2 n - 1}}$$

$$= \frac{1 - r^{\log_2 n}}{1 - r}$$

$$= \frac{1 - 2^{\log_2 n}}{1 - 2}$$

$$= 2^{\log_2 n} - 1$$

*

$$= \log_2 n \cdot 2^{\log_2 n} - (2^{\log_2 n} - 1)$$

$$= n \log_2 n - (n - 1)$$

gcd : after two iterations of

gcd(a, b) in the worst case

I'm down to

gcd($\frac{a}{2}$, stuff)

$$C(n) = 1 \cdot C\left(\frac{n}{2}\right) + 2$$

$$C(1) = 0$$

$$c = 1$$

$$g(n) = 2$$

$$a = 0$$

$$C(n) = 0 + \sum_{i=1}^{\log n} 1^{\log n - i} \cdot 2$$
$$= \sum_{i=1}^{\log n} 2$$

$$C(n) = 2 \log n$$