

Algorithms

Copyright ©2006 S. Dasgupta, C. H. Papadimitriou, and U. V. Vazirani

July 18, 2006

Contents

Preface	9
0 Prologue	11
0.1 Books and algorithms	11
0.2 Enter Fibonacci	12
0.3 Big- O notation	15
Exercises	18
1 Algorithms with numbers	21
1.1 Basic arithmetic	21
1.2 Modular arithmetic	25
1.3 Primality testing	33
1.4 Cryptography	39
1.5 Universal hashing	43
Exercises	48
Randomized algorithms: a virtual chapter	39
2 Divide-and-conquer algorithms	55
2.1 Multiplication	55
2.2 Recurrence relations	58
2.3 Mergesort	60
2.4 Medians	64
2.5 Matrix multiplication	66
2.6 The fast Fourier transform	68
Exercises	83
3 Decompositions of graphs	91
3.1 Why graphs?	91
3.2 Depth-first search in undirected graphs	93
3.3 Depth-first search in directed graphs	98
3.4 Strongly connected components	101
Exercises	106

4	Paths in graphs	115
4.1	Distances	115
4.2	Breadth-first search	116
4.3	Lengths on edges	118
4.4	Dijkstra's algorithm	119
4.5	Priority queue implementations	126
4.6	Shortest paths in the presence of negative edges	128
4.7	Shortest paths in dags	130
	Exercises	132
5	Greedy algorithms	139
5.1	Minimum spanning trees	139
5.2	Huffman encoding	153
5.3	Horn formulas	157
5.4	Set cover	158
	Exercises	161
6	Dynamic programming	169
6.1	Shortest paths in dags, revisited	169
6.2	Longest increasing subsequences	170
6.3	Edit distance	174
6.4	Knapsack	181
6.5	Chain matrix multiplication	184
6.6	Shortest paths	186
6.7	Independent sets in trees	189
	Exercises	191
7	Linear programming and reductions	201
7.1	An introduction to linear programming	201
7.2	Flows in networks	211
7.3	Bipartite matching	219
7.4	Duality	220
7.5	Zero-sum games	224
7.6	The simplex algorithm	227
7.7	Postscript: circuit evaluation	236
	Exercises	239
8	NP-complete problems	247
8.1	Search problems	247
8.2	NP-complete problems	257
8.3	The reductions	262
	Exercises	278

9 Coping with NP-completeness	283
9.1 Intelligent exhaustive search	284
9.2 Approximation algorithms	290
9.3 Local search heuristics	297
Exercises	306
10 Quantum algorithms	311
10.1 Qubits, superposition, and measurement	311
10.2 The plan	315
10.3 The quantum Fourier transform	316
10.4 Periodicity	318
10.5 Quantum circuits	321
10.6 Factoring as periodicity	324
10.7 The quantum algorithm for factoring	326
Exercises	329
Historical notes and further reading	331
Index	333

List of boxes

Bases and logs	21
Two's complement	27
Is your social security number a prime?	33
Hey, that was group theory!	36
Carmichael numbers	37
Randomized algorithms: a virtual chapter	39
An application of number theory?	40
Binary search	60
An $n \log n$ lower bound for sorting	62
The Unix <code>sort</code> command	66
Why multiply polynomials?	68
The slow spread of a fast algorithm	82
How big is your graph?	93
Crawling fast	105
Which heap is best?	125
Trees	140
A randomized algorithm for minimum cut	150
Entropy	155
Recursion? No, thanks.	173
Programming?	173
Common subproblems	177
Of mice and men	179
Memoization	183
On time and memory	189
A magic trick called duality	205
Reductions	209
Matrix-vector notation	211
Visualizing duality	222
Gaussian elimination	234

Linear programming in polynomial time	236
The story of Sissa and Moore	247
Why P and NP ?	258
The two ways to use reductions	259
Unsolvable problems	276
Entanglement	314
The Fourier transform of a periodic vector	320
Setting up a periodic superposition	325
Quantum physics meets computation	327

Preface

This book evolved over the past ten years from a set of lecture notes developed while teaching the undergraduate Algorithms course at Berkeley and U.C. San Diego. Our way of teaching this course evolved tremendously over these years in a number of directions, partly to address our students' background (undeveloped formal skills outside of programming), and partly to reflect the maturing of the field in general, as we have come to see it. The notes increasingly crystallized into a narrative, and we progressively structured the course to emphasize the “story line” implicit in the progression of the material. As a result, the topics were carefully selected and clustered. No attempt was made to be encyclopedic, and this freed us to include topics traditionally de-emphasized or omitted from most Algorithms books.

Playing on the strengths of our students (shared by most of today's undergraduates in Computer Science), instead of dwelling on formal proofs we distilled in each case the crisp mathematical idea that makes the algorithm work. In other words, we emphasized rigor over formalism. We found that our students were much more receptive to mathematical rigor of this form. It is this progression of crisp ideas that helps weave the story.

Once you think about Algorithms in this way, it makes sense to start at the historical beginning of it all, where, in addition, the characters are familiar and the contrasts dramatic: numbers, primality, and factoring. This is the subject of Part I of the book, which also includes the RSA cryptosystem, and divide-and-conquer algorithms for integer multiplication, sorting and median finding, as well as the fast Fourier transform. There are three other parts: Part II, the most traditional section of the book, concentrates on data structures and graphs; the contrast here is between the intricate structure of the underlying problems and the short and crisp pieces of pseudocode that solve them. Instructors wishing to teach a more traditional course can simply start with Part II, which is self-contained (following the prologue), and then cover Part I as required. In Parts I and II we introduced certain techniques (such as greedy and divide-and-conquer) which work for special kinds of problems; Part III deals with the “sledgehammers” of the trade, techniques that are powerful and general: dynamic programming (a novel approach helps clarify this traditional stumbling block for students) and linear programming (a clean and intuitive treatment of the simplex algorithm, duality, and reductions to the basic problem). The final Part IV is about ways of dealing with hard problems: NP-completeness, various heuristics, as well as quantum algorithms, perhaps the most advanced and modern topic. As it happens, we end the story exactly where we started it, with Shor's quantum algorithm for factoring.

The book includes three additional undercurrents, in the form of three series of separate

“boxes,” strengthening the narrative (and addressing variations in the needs and interests of the students) while keeping the flow intact: pieces that provide historical context; descriptions of how the explained algorithms are used in practice (with emphasis on internet applications); and excursions for the mathematically sophisticated.

Chapter 0

Prologue

Look around you. Computers and networks are everywhere, enabling an intricate web of complex human activities: education, commerce, entertainment, research, manufacturing, health management, human communication, even war. Of the two main technological underpinnings of this amazing proliferation, one is obvious: the breathtaking pace with which advances in microelectronics and chip design have been bringing us faster and faster hardware.

This book tells the story of the other intellectual enterprise that is crucially fueling the computer revolution: *efficient algorithms*. It is a fascinating story.

Gather 'round and listen close.

0.1 Books and algorithms

Two ideas changed the world. In 1448 in the German city of Mainz a goldsmith named Johann Gutenberg discovered a way to print books by putting together movable metallic pieces. Literacy spread, the Dark Ages ended, the human intellect was liberated, science and technology triumphed, the Industrial Revolution happened. Many historians say we owe all this to typography. Imagine a world in which only an elite could read these lines! But others insist that the key development was not typography, but *algorithms*.

Today we are so used to writing numbers in decimal, that it is easy to forget that Gutenberg would write the number 1448 as MCDXLVIII. How do you add two Roman numerals? What is MCDXLVIII + DCCCXII? (And just try to think about multiplying them.) Even a clever man like Gutenberg probably only knew how to add and subtract small numbers using his fingers; for anything more complicated he had to consult an abacus specialist.

The decimal system, invented in India around AD 600, was a revolution in quantitative reasoning: using only 10 symbols, even very large numbers could be written down compactly, and arithmetic could be done efficiently on them by following elementary steps. Nonetheless these ideas took a long time to spread, hindered by traditional barriers of language, distance, and ignorance. The most influential medium of transmission turned out to be a textbook, written in Arabic in the ninth century by a man who lived in Baghdad. Al Khwarizmi laid out the basic methods for adding, multiplying, and dividing numbers—even extracting square roots and calculating digits of π . These procedures were precise, unambiguous, mechanical,

efficient, correct—in short, they were *algorithms*, a term coined to honor the wise man after the decimal system was finally adopted in Europe, many centuries later.

Since then, this decimal positional system and its numerical algorithms have played an enormous role in Western civilization. They enabled science and technology; they accelerated industry and commerce. And when, much later, the computer was finally designed, it explicitly embodied the positional system in its bits and words and arithmetic unit. Scientists everywhere then got busy developing more and more complex algorithms for all kinds of problems and inventing novel applications—ultimately changing the world.

0.2 Enter Fibonacci

Al Khwarizmi’s work could not have gained a foothold in the West were it not for the efforts of one man: the 15th century Italian mathematician Leonardo Fibonacci, who saw the potential of the positional system and worked hard to develop it further and propagandize it.

But today Fibonacci is most widely known for his famous sequence of numbers

$$0, 1, 1, 2, 3, 5, 8, 13, 21, 34, \dots,$$

each the sum of its two immediate predecessors. More formally, the Fibonacci numbers F_n are generated by the simple rule

$$F_n = \begin{cases} F_{n-1} + F_{n-2} & \text{if } n > 1 \\ 1 & \text{if } n = 1 \\ 0 & \text{if } n = 0. \end{cases}$$

No other sequence of numbers has been studied as extensively, or applied to more fields: biology, demography, art, architecture, music, to name just a few. And, together with the powers of 2, it is computer science’s favorite sequence.

In fact, the Fibonacci numbers grow *almost* as fast as the powers of 2: for example, F_{30} is over a million, and F_{100} is already 21 digits long! In general, $F_n \approx 2^{0.694n}$ (see Exercise 0.3).

But what is the precise value of F_{100} , or of F_{200} ? Fibonacci himself would surely have wanted to know such things. To answer, we need an algorithm for computing the n th Fibonacci number.

An exponential algorithm

One idea is to slavishly implement the recursive definition of F_n . Here is the resulting algorithm, in the “pseudocode” notation used throughout this book:

```
function fib1(n)
if n = 0: return 0
if n = 1: return 1
return fib1(n - 1) + fib1(n - 2)
```

Whenever we have an algorithm, there are three questions we always ask about it:

1. Is it correct?
2. How much time does it take, as a function of n ?
3. And can we do better?

The first question is moot here, as this algorithm is precisely Fibonacci's definition of F_n . But the second demands an answer. Let $T(n)$ be the number of *computer steps* needed to compute `fib1`(n); what can we say about this function? For starters, if n is less than 2, the procedure halts almost immediately, after just a couple of steps. Therefore,

$$T(n) \leq 2 \text{ for } n \leq 1.$$

For larger values of n , there are two recursive invocations of `fib1`, taking time $T(n-1)$ and $T(n-2)$, respectively, plus three computer steps (checks on the value of n and a final addition). Therefore,

$$T(n) = T(n-1) + T(n-2) + 3 \text{ for } n > 1.$$

Compare this to the recurrence relation for F_n : we immediately see that $T(n) \geq F_n$.

This is very bad news: the running time of the algorithm grows as fast as the Fibonacci numbers! $T(n)$ is *exponential in n* , which implies that the algorithm is impractically slow except for very small values of n .

Let's be a little more concrete about just how bad exponential time is. To compute F_{200} , the `fib1` algorithm executes $T(200) \geq F_{200} \geq 2^{138}$ elementary computer steps. How long this actually takes depends, of course, on the computer used. At this time, the fastest computer in the world is the NEC Earth Simulator, which clocks 40 trillion steps per second. Even on this machine, `fib1`(200) would take at least 2^{92} seconds. This means that, if we start the computation today, it would still be going long after the sun turns into a red giant star.

But technology is rapidly improving—computer speeds have been doubling roughly every 18 months, a phenomenon sometimes called *Moore's law*. With this extraordinary growth, perhaps `fib1` will run a lot faster on next year's machines. Let's see—the running time of `fib1`(n) is proportional to $2^{0.694n} \approx (1.6)^n$, so it takes 1.6 times longer to compute F_{n+1} than F_n . And under Moore's law, computers get roughly 1.6 times faster each year. So if we can reasonably compute F_{100} with this year's technology, then next year we will manage F_{101} . And the year after, F_{102} . And so on: just one more Fibonacci number every year! Such is the curse of exponential time.

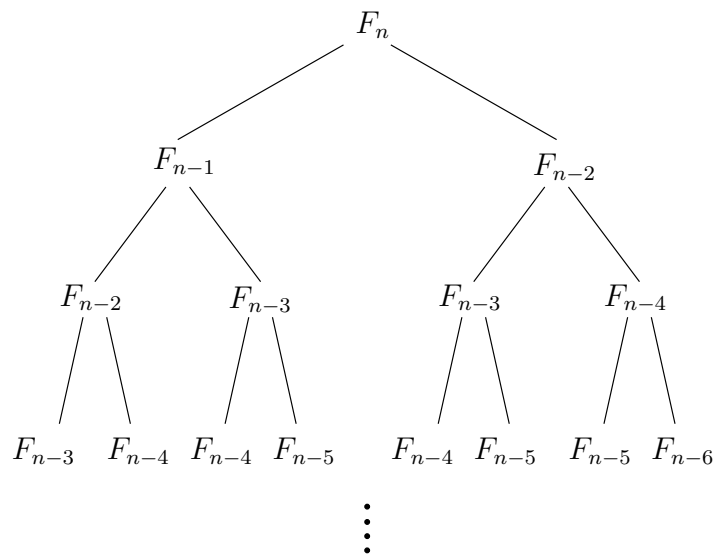
In short, our naive recursive algorithm is correct but hopelessly inefficient. *Can we do better?*

A polynomial algorithm

Let's try to understand why `fib1` is so slow. Figure 0.1 shows the cascade of recursive invocations triggered by a single call to `fib1`(n). Notice that many computations are repeated!

A more sensible scheme would store the intermediate results—the values F_0, F_1, \dots, F_{n-1} —as soon as they become known.

Figure 0.1 The proliferation of recursive calls in `fib1`.



```

function fib2(n)
  if n = 0 return 0
  create an array f[0...n]
  f[0] = 0, f[1] = 1
  for i = 2...n:
    f[i] = f[i-1] + f[i-2]
  return f[n]

```

As with `fib1`, the correctness of this algorithm is self-evident because it directly uses the definition of F_n . How long does it take? The inner loop consists of a single computer step and is executed $n - 1$ times. Therefore the number of computer steps used by `fib2` is *linear in n* . From exponential we are down to *polynomial*, a huge breakthrough in running time. It is now perfectly reasonable to compute F_{200} or even $F_{200,000}$.¹

As we will see repeatedly throughout this book, the right algorithm makes all the difference.

More careful analysis

In our discussion so far, we have been counting the number of *basic computer steps* executed by each algorithm and thinking of these basic steps as taking a constant amount of time. This is a very useful simplification. After all, a processor's instruction set has a variety of basic primitives—branching, storing to memory, comparing numbers, simple arithmetic, and

¹To better appreciate the importance of this dichotomy between exponential and polynomial algorithms, the reader may want to peek ahead to *the story of Sissa and Moore*, in Chapter 8.

so on—and rather than distinguishing between these elementary operations, it is far more convenient to lump them together into one category.

But looking back at our treatment of Fibonacci algorithms, we have been too liberal with what we consider a basic step. It is reasonable to treat addition as a single computer step if small numbers are being added, 32-bit numbers say. But the n th Fibonacci number is about $0.694n$ bits long, and this can far exceed 32 as n grows. Arithmetic operations on arbitrarily large numbers cannot possibly be performed in a single, constant-time step. We need to audit our earlier running time estimates and make them more honest.

We will see in Chapter 1 that the addition of two n -bit numbers takes time roughly proportional to n ; this is not too hard to understand if you think back to the grade-school procedure for addition, which works on one digit at a time. Thus `fib1`, which performs about F_n additions, actually uses a number of *basic steps* roughly proportional to nF_n . Likewise, the number of steps taken by `fib2` is proportional to n^2 , still polynomial in n and therefore exponentially superior to `fib1`. This correction to the running time analysis does not diminish our breakthrough.

But can we do even better than fib2? Indeed we can: see Exercise 0.4.

0.3 Big- O notation

We’ve just seen how sloppiness in the analysis of running times can lead to an unacceptable level of inaccuracy in the result. But the opposite danger is also present: it is possible to be *too* precise. An insightful analysis is based on the right simplifications.

Expressing running time in terms of *basic computer steps* is already a simplification. After all, the time taken by one such step depends crucially on the particular processor and even on details such as caching strategy (as a result of which the running time can differ subtly from one execution to the next). Accounting for these architecture-specific minutiae is a nightmarishly complex task and yields a result that does not generalize from one computer to the next. It therefore makes more sense to seek an uncluttered, machine-independent characterization of an algorithm’s efficiency. To this end, we will always express running time by counting the number of basic computer steps, as a function of the size of the input.

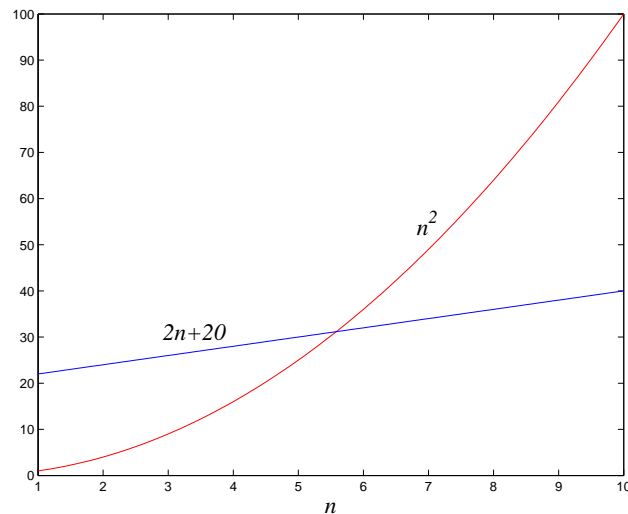
And this simplification leads to another. Instead of reporting that an algorithm takes, say, $5n^3 + 4n + 3$ steps on an input of size n , it is much simpler to leave out lower-order terms such as $4n$ and 3 (which become insignificant as n grows), and even the detail of the coefficient 5 in the leading term (computers will be five times faster in a few years anyway), and just say that the algorithm takes time $O(n^3)$ (pronounced “big oh of n^3 ”).

It is time to define this notation precisely. In what follows, think of $f(n)$ and $g(n)$ as the running times of two algorithms on inputs of size n .

Let $f(n)$ and $g(n)$ be functions from positive integers to positive reals. We say $f = O(g)$ (which means that “ f grows no faster than g ”) if *there is a constant $c > 0$ such that $f(n) \leq c \cdot g(n)$.*

Saying $f = O(g)$ is a very loose analog of “ $f \leq g$.” It differs from the usual notion of \leq because of the constant c , so that for instance $10n = O(n)$. This constant also allows us to

Figure 0.2 Which running time is better?



disregard what happens for small values of n . For example, suppose we are choosing between two algorithms for a particular computational task. One takes $f_1(n) = n^2$ steps, while the other takes $f_2(n) = 2n + 20$ steps (Figure 0.2). Which is better? Well, this depends on the value of n . For $n \leq 5$, f_1 is smaller; thereafter, f_2 is the clear winner. In this case, f_2 scales much better as n grows, and therefore it is superior.

This superiority is captured by the big- O notation: $f_2 = O(f_1)$, because

$$\frac{f_2(n)}{f_1(n)} = \frac{2n + 20}{n^2} \leq 22$$

for all n ; on the other hand, $f_1 \neq O(f_2)$, since the ratio $f_1(n)/f_2(n) = n^2/(2n + 20)$ can get arbitrarily large, and so no constant c will make the definition work.

Now another algorithm comes along, one that uses $f_3(n) = n + 1$ steps. Is this better than f_2 ? Certainly, but only by a constant factor. The discrepancy between f_2 and f_3 is tiny compared to the huge gap between f_1 and f_2 . In order to stay focused on the big picture, we treat functions as equivalent if they differ only by multiplicative constants.

Returning to the definition of big- O , we see that $f_2 = O(f_3)$:

$$\frac{f_2(n)}{f_3(n)} = \frac{2n + 20}{n + 1} \leq 20,$$

and of course $f_3 = O(f_2)$, this time with $c = 1$.

Just as $O(\cdot)$ is an analog of \leq , we can also define analogs of \geq and $=$ as follows:

$$\begin{aligned} f = \Omega(g) &\text{ means } g = O(f) \\ f = \Theta(g) &\text{ means } f = O(g) \text{ and } f = \Omega(g). \end{aligned}$$

In the preceding example, $f_2 = \Theta(f_3)$ and $f_1 = \Omega(f_3)$.

Big- O notation lets us focus on the big picture. When faced with a complicated function like $3n^2 + 4n + 5$, we just replace it with $O(f(n))$, where $f(n)$ is as simple as possible. In this particular example we'd use $O(n^2)$, because the quadratic portion of the sum dominates the rest. Here are some commonsense rules that help simplify functions by omitting dominated terms:

1. Multiplicative constants can be omitted: $14n^2$ becomes n^2 .
2. n^a dominates n^b if $a > b$: for instance, n^2 dominates n .
3. Any exponential dominates any polynomial: 3^n dominates n^5 (it even dominates 2^n).
4. Likewise, any polynomial dominates any logarithm: n dominates $(\log n)^3$. This also means, for example, that n^2 dominates $n \log n$.

Don't misunderstand this cavalier attitude toward constants. Programmers and algorithm developers are *very* interested in constants and would gladly stay up nights in order to make an algorithm run faster by a factor of 2. But understanding algorithms at the level of this book would be impossible without the simplicity afforded by big- O notation.