

# Algorithms

Department of Computer Science  
University of Illinois at Urbana-Champaign

Instructor: Jeff Erickson

## Teaching Assistants:

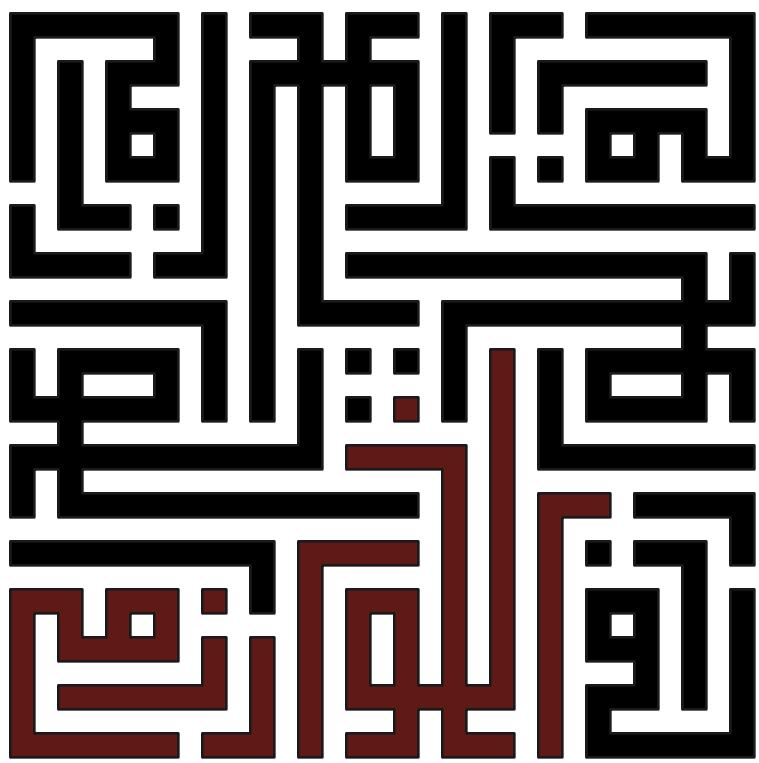
- **Spring 1999:** Mitch Harris and Shripad Thite
- **Summer 1999 (IMCS):** Mitch Harris
- **Summer 2000 (IMCS):** Mitch Harris
- **Fall 2000:** Chris Neihengen, Ekta Manaktala, and Nick Hurlburt
- **Spring 2001:** Brian Ensink, Chris Neihengen, and Nick Hurlburt
- **Summer 2001 (I2CS):** Asha Seetharam and Dan Bullok
- **Fall 2002:** Erin Wolf, Gio Kao, Kevin Small, Michael Bond, Rishi Talreja, Rob McCann, and Yasutaka Furakawa
- **Spring 2004:** Dan Cranston, Johnathon Fischer, Kevin Milans, and Lan Chen
- **Fall 2005:** Erin Chambers, Igor Gammer, and Aditya Ramani
- **Fall 2006:** Dan Cranston, Nitish Korula, and Kevin Milans
- **Spring 2007:** Kevin Milans
- **Fall 2008:** Reza Zamani-Nasab
- **Spring 2009:** Alina Ene, Ben Moseley, and Amir Nayyeri
- **Spring 2010:** David Morrison, Kyle Fox, and Rachit Agarwal
- **Fall 2010:** Alina Ene

© Copyright 1999–2011 Jeff Erickson. Last update July 8, 2011.

This work may be freely copied and distributed, either electronically or on paper.  
It may not be sold for more than the actual cost of reproduction, storage, or transmittal.

This work is licensed under a Creative Commons Attribution-NonCommercial-Share Alike 3.0 United States License.  
For license details, see <http://creativecommons.org/licenses/by-nc-sa/3.0/us/>.

For the most recent edition of this work, see <http://www.cs.illinois.edu/~jeffe/teaching/algorithms/>.



*Shall I tell you, my friend, how you will come to understand it?  
Go and write a book on it.*

— Henry Home, Lord Kames (1696–1782), to Sir Gilbert Elliot

*You know, I could write a book.  
And this book would be thick enough to stun an ox.*

— Laurie Anderson, “Let X=X”, *Big Science* (1982)

*I’m writing a book. I’ve got the page numbers done,  
so now I just have to fill in the rest.*

— Stephen Wright

## About These Notes

This course packet includes lecture notes, homework questions, and exam questions from algorithms courses I taught at the University of Illinois at Urbana-Champaign in Spring 1999, Fall 2000, Spring 2001, Fall 2002, Spring 2004, Fall 2005, Fall 2006, Spring 2007, Fall 2008, Spring 2009, Spring 2010, and Fall 2010. These lecture notes and my videotaped lectures were also offered over the web in Summer 1999, Summer 2000, Summer 2001, Fall 2002, and Fall 2005 as part of the UIUC computer science department’s online master’s program. Lecture notes were posted to the course web site a few days (on average) after each lecture. Homeworks, exams, and solutions were also distributed over the web.

Most (but not all) of the exercises at the end of each lecture note have been used at least once in a homework assignment, discussion section, or exam. You can also find a near-complete collection of homeworks and exams from past semesters of my class online at <http://www.cs.illinois.edu/~jeffe/teaching/algorithms/>. A large fraction of these exercises were contributed by some amazing teaching assistants:

Aditya Ramani, Alina Ene, Amir Nayyeri, Asha Seetharam, Ben Moseley, Brian Ensink, Chris Neihengen, Dan Bullok, Dan Cranston, David Morrison, Johnathon Fischer, Ekta Manaktala, Erin Wolf Chambers, Igor Gammer, Gio Kao, Kevin Milans, Kevin Small, Kyle Fox, Lan Chen, Michael Bond, Mitch Harris, Nick Hurlburt, Nitish Korula, Rachit Agarwal, Reza Zamani-Nasab, Rishi Talreja, Rob McCann, Shripad Thite, and Yasu Furakawa.

\*Stars indicate more challenging problems; many of these appeared on qualifying exams for the algorithms PhD students at UIUC. A small number of *really* hard problems are marked with a ★larger star; one or two *open* problems are indicated by ★enormous stars.

Please do not ask me for solutions to the exercises. If you’re a student, seeing the solution will rob you of the experience of solving the problem yourself, which is the only way to learn the material. If you’re an instructor, you shouldn’t assign problems that you can’t solve yourself! (I do not always follow my own advice; some of these problems have serious bugs.)

## Acknowledgments

The lecture notes and exercises draw heavily on the creativity, wisdom, and experience of thousands of algorithms students, teachers, and researchers. In particular, I am immensely grateful to the almost 1400 Illinois students who have used these notes as a primary reference, offered useful (if sometimes painful) criticism, and suffered through some truly awful first drafts. I'm also grateful for the contributions and feedback from teaching assistants, all listed above.

Naturally, these notes owe a great deal to the people who taught me this algorithms stuff in the first place: Bob Bixby and Michael Perlman at Rice; David Eppstein, Dan Hirshberg, and George Lueker at UC Irvine; and Abhiram Ranade, Dick Karp, Manuel Blum, Mike Luby, and Raimund Seidel at UC Berkeley. I've also been helped tremendously by many discussions with faculty colleagues at UIUC—Edgar Ramos, Herbert Edelsbrunner, Jason Zych, Lenny Pitt, Mahesh Viswanathan, Margaret Fleck, Shang-Hua Teng, Steve LaValle, and especially Chandra Chekuri, Ed Reingold, and Sariel Har-Peled. I stole the first iteration of the overall course structure, and the idea to write up my own lecture notes, from Herbert Edelsbrunner.

Finally, “Johnny’s” multi-colored crayon homework was found under the TA office door among the other Fall 2000 Homework 1 submissions. The square Kufi rendition of the name “al-Khwārizmī” on the back of the cover page is mine.

## Additional References

I strongly encourage my students (and other readers) *not* to restrict themselves to a single textual reference. Authors and readers bring their own perspectives to the material; no instructor ‘clicks’ with every student, or even every very strong student. Finding the author that most effectively gets their intuition into *your* head take some effort, but that effort pays off handsomely in the long run.

The following references have been particularly valuable sources of inspiration, intuition, examples, and problems. This list is incomplete!

- Alfred V. Aho, John E. Hopcroft, and Jeffrey D. Ullman. *The Design and Analysis of Computer Algorithms*. Addison-Wesley, 1974. (I used this textbook as an undergrad at Rice, and again as a masters student at UC Irvine.)
- Thomas Cormen, Charles Leiserson, Ron Rivest, and Cliff Stein. *Introduction to Algorithms*, third edition. MIT Press/McGraw-Hill, 2009. (The second edition was my recommended textbook until 2005. I also used the first edition as a teaching assistant at Berkeley.)
- Sanjoy Dasgupta, Christos H. Papadimitriou, and Umesh V. Vazirani. *Algorithms*. McGraw-Hill, 2006. (This is the current recommended textbook for my undergraduate classes.)
- Jeff Edmonds. *How to Think about Algorithms*. Cambridge University Press, 2008.
- Michael R. Garey and David S. Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman, 1979.
- Michael T. Goodrich and Roberto Tamassia. *Algorithm Design: Foundations, Analysis, and Internet Examples*. John Wiley & Sons, 2002.
- Jon Kleinberg and Éva Tardos. *Algorithm Design*. Addison-Wesley, 2005. (This is the current recommended textbook for my graduate algorithms classes.)
- Donald Knuth. *The Art of Computer Programming*, volumes 1–3. Addison-Wesley, 1997. (My parents gave me these for Christmas when I was 14. I didn't actually read them until much later.)

- Udi Manber. *Introduction to Algorithms: A Creative Approach*. Addison-Wesley, 1989. (I used this textbook as a teaching assistant at Berkeley.)
- Rajeev Motwani and Prabhakar Raghavan. *Randomized Algorithms*. Cambridge University Press, 1995.
- Ian Parberry. *Problems on Algorithms*. Prentice-Hall, 1995. (This book is out of print, but it can be downloaded for ‘free’ from <http://www.eng.unt.edu/ian/books/free/license.html>.)
- Alexander Schrijver. *Combinatorial Optimization: Polyhedra and Efficiency*. Springer, 2003.
- Robert Sedgewick. *Algorithms*. Addison-Wesley, 1988. (This book and its sequels have by far the best algorithm *illustrations* I’ve seen anywhere.)
- Robert Endre Tarjan. *Data Structures and Network Algorithms*. SIAM, 1983.
- Robert J. Vanderbei. *Linear Programming: Foundations and Extensions*. Springer, 2001.
- Class notes from my own algorithms classes at Berkeley, especially those taught by Dick Karp and Raimund Seidel.
- Lecture notes, slides, homeworks, exams, and video lectures posted by innumerable colleagues around the world.
- The Source of All Knowledge (Google) and The Source of All Lies (Wikipedia).

## Prerequisites

For the most part, these notes assume the reader has mastered the material covered in the first two years of a strong undergraduate computer science curriculum, and has the intellectual maturity to recognize and repair any remaining gaps in their mastery. (Mastery is not the same thing as ‘exposure’ or ‘a good grade’; this is why I start every semester with Homework Zero.) Specific prerequisites include:

- Proof techniques: direct proof, indirect proof, proof by contradiction, combinatorial proof, and induction (including its “strong”, “structural”, and “recursive” forms). Lecture 0 requires induction, and whenever Lecture  $n - 1$  requires induction, so does Lecture  $n$ .
- Discrete mathematics: High-school algebra, naive set theory, Boolean algebra, first-order predicate logic, sets, functions, relations, modular arithmetic, recursive definitions, trees (as abstract objects, not data structures), graphs.
- Elementary discrete probability: uniform vs non-uniform probability distributions, expectation, linearity of expectation, independence.
- Iterative programming concepts: variables, conditionals, iteration, subroutines, indirection (addresses/pointers/references), recursion. Programming experience in any language that supports pointers and recursion is a plus.
- Fundamental data structures: arrays, linked lists, binary search trees, at least one balanced search tree (such as AVL trees, red-black trees, B-trees, skip lists, splay trees, or treaps), binary heaps.
- Fundamental abstract data types: dictionaries, stacks, queues, priority queues; the difference between this list and the previous list.
- Fundamental algorithms: elementary arithmetic, linear search, binary search, sorting (selection, insertion, merge-, heap-, quick-, radix, anything but bubble-), pre-/post-/inorder tree traversal.

- Basic algorithm analysis: Asymptotic notation ( $o$ ,  $O$ ,  $\Theta$ ,  $\Omega$ ,  $\omega$ ), translating loops into sums and recursive calls into recurrences, evaluating simple sums and recurrences.
- Mathematical maturity: facility with abstraction, formal (especially recursive) definitions, and (especially inductive) proofs; following mathematical arguments; recognizing syntactic, semantic, and/or logical nonsense; writing the former rather than the latter.

Some of this prerequisite material is covered briefly in these notes, but more as a reminder than a good introduction.

### Caveat Lector!

With few exceptions, each of these notes contains far too much material to cover in a single lecture. In a typical 75-minute lecture, I tend to cover 4 to 5 pages of material—a bit more if I'm lecturing to graduate students than to undergraduates. Your mileage may vary! (Arguably, that means that as I continue to add material, the label “lecture notes” becomes less and less accurate.)

Despite several rounds of revision, these notes still contain lots of mistakes, errors, bugs, gaffes, omissions, snafus, kludges, typos, mathos, grammaros, thinkos, brain farts, nonsense, garbage, cruft, junk, and outright lies, ***all of which are entirely Steve Skiena's fault.*** I revise and update these notes every time I teach the course, so please let me know if you find a bug. (Steve is unlikely to care.)

Whenever I teach the algorithms class, I award extra credit points to the first student to post an explanation and correction of any error in the lecture notes to the course newsgroup. Obviously, the number of extra credit points depends on the severity of the error and the quality of the correction. If I'm not teaching the course, encourage your instructor to set up a similar extra-credit scheme, and forward the bug reports to **Steve** me!

Of course, any other feedback is also welcome!

Enjoy!

— Jeff

*It is traditional for the author to magnanimously accept the blame for whatever deficiencies remain. I don't. Any errors, deficiencies, or problems in this book are somebody else's fault, but I would appreciate knowing about them so as to determine who is to blame.*

— Steven S. Skiena, *The Algorithm Design Manual* (1997)

*We should explain, before proceeding, that it is not our object to consider this program with reference to the actual arrangement of the data on the Variables of the engine, but simply as an abstract question of the nature and number of the operations required to be performed during its complete solution.*

— Ada Augusta Byron King, Countess of Lovelace, translator's notes for Luigi F. Menabrea,  
"Sketch of the Analytical Engine invented by Charles Babbage, Esq." (1843)

*You are right to demand that an artist engage his work consciously, but you confuse two different things: solving the problem and correctly posing the question.*

— Anton Chekhov, in a letter to A. S. Suvorin (October 27, 1888)

*The more we reduce ourselves to machines in the lower things,  
the more force we shall set free to use in the higher.*

— Anna C. Brackett, *The Technique of Rest* (1892)

*The moment a man begins to talk about technique  
that's proof that he is fresh out of ideas.*

— Raymond Chandler

## 0 Introduction

### 0.1 What is an algorithm?

An algorithm is an explicit, precise, unambiguous, mechanically-executable sequence of elementary instructions. For example, here is an algorithm for singing that annoying song '99 Bottles of Beer on the Wall', for arbitrary values of 99:

```
BOTTLESOFBEER( $n$ ):  
  For  $i \leftarrow n$  down to 1  
    Sing " $i$  bottles of beer on the wall,  $i$  bottles of beer,"  
    Sing "Take one down, pass it around,  $i - 1$  bottles of beer on the wall."  
  
  Sing "No bottles of beer on the wall, no bottles of beer,"  
  Sing "Go to the store, buy some more,  $n$  bottles of beer on the wall."
```

The word 'algorithm' does *not* derive, as algorithmophobic classicists might guess, from the Greek root *algos* (ἀλγός), meaning 'pain'. Rather, it is a corruption of the name of the 9th century Persian mathematician Abū 'Abd Allāh Muḥammad ibn Mūsā al-Khwārizmī.<sup>1</sup> Al-Khwārizmī is perhaps best known as the writer of the treatise *Al-Kitāb al-mukhtaṣar fīhiṣāb al-ğabr wa'l-muqābala*<sup>2</sup>, from which the modern word *algebra* derives. In another treatise, al-Khwārizmī popularized the modern decimal system for writing and manipulating numbers—in particular, the use of a small circle or *sifr* to represent a missing quantity—which had originated in India several centuries earlier. This system later became known in Europe as *algorism*. Thanks to the efforts of the medieval Italian mathematician Leonardo of Pisa, better known as Fibonacci, algorism began to replace the abacus as the preferred system of commercial calculation<sup>3</sup> in Europe in the late 12th century, although cyphers became truly ubiquitous in Western Europe only after the French revolution 600 years later. The more modern word *algorithm* is a false

<sup>1</sup>'Mohammad, father of Adbdulla, son of Moses, the Kwārizmian'. Kwārizm is an ancient city, now called Khiva, in the Khorezm Province of Uzbekistan.

<sup>2</sup>The Compendious Book on Calculation by Completion and Balancing'.

<sup>3</sup>from the Latin word *calculus*, meaning literally 'small rock', referring to the stones on a counting board, or abacus

cognate with the Greek word *arithmos* (ἀριθμός), meaning ‘number’ (and perhaps the aforementioned ἀλγός). Thus, until very recently, the word *algorithm* referred exclusively to pencil-and-paper methods for numerical calculations. People trained in the reliable execution of these methods were called—you guessed it—*computers*.

## 0.2 A Few Simple Examples

### Multiplication by compass and straightedge

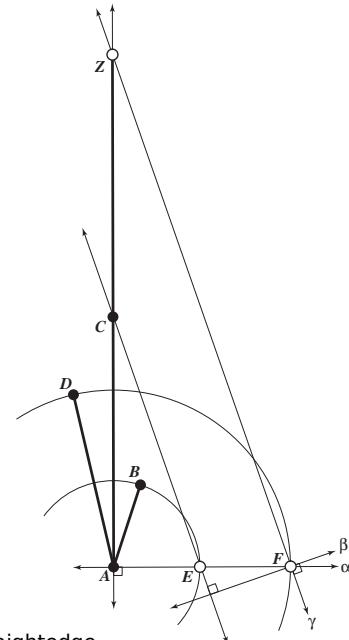
Although they have only been an object of formal study for a few decades, algorithms have been with us since the dawn of civilization, for centuries before Al-Khwārizmī and Fibonacci popularized the cypher. Here is an algorithm, popularized (but almost certainly not discovered) by Euclid about 2500 years ago, for multiplying or dividing numbers using a ruler and compass. The Greek geometers represented numbers using line segments of the appropriate length. In the pseudo-code below,  $\text{CIRCLE}(p, q)$  represents the circle centered at a point  $p$  and passing through another point  $q$ . Hopefully the other instructions are obvious.<sup>4</sup>

```

⟨⟨Construct the line perpendicular to  $\ell$  and passing through  $P$ .⟩⟩
RIGHTANGLE( $\ell, P$ ):
    Choose a point  $A \in \ell$ 
     $A, B \leftarrow \text{INTERSECT}(\text{CIRCLE}(P, A), \ell)$ 
     $C, D \leftarrow \text{INTERSECT}(\text{CIRCLE}(A, B), \text{CIRCLE}(B, A))$ 
    return LINE( $C, D$ )
⟨⟨Construct a point  $Z$  such that  $|AZ| = |AC||AD|/|AB|$ .⟩⟩
MULTIPLYOrDIVIDE( $A, B, C, D$ ):
     $\alpha \leftarrow \text{RIGHTANGLE}(\text{LINE}(A, C), A)$ 
     $E \leftarrow \text{INTERSECT}(\text{CIRCLE}(A, B), \alpha)$ 
     $F \leftarrow \text{INTERSECT}(\text{CIRCLE}(A, D), \alpha)$ 
     $\beta \leftarrow \text{RIGHTANGLE}(\text{LINE}(E, C), F)$ 
     $\gamma \leftarrow \text{RIGHTANGLE}(\beta, F)$ 
    return  $\text{INTERSECT}(\gamma, \text{LINE}(A, C))$ 

```

Multiplying or dividing using a compass and straightedge.



This algorithm breaks down the difficult task of multiplication into a series of simple primitive operations: drawing a line between two points, drawing a circle with a given center and boundary point, and so on. These primitive steps are quite non-trivial to execute on a modern digital computer, but this algorithm wasn’t designed for a digital computer; it was designed for the Platonic Ideal Classical Greek Mathematician, wielding the Platonic Ideal Compass and the Platonic Ideal Straightedge. In this example, Euclid first defines a new primitive operation, constructing a right angle, by (as modern programmers would put it) writing a subroutine.

<sup>4</sup>Euclid and his students almost certainly drew their constructions on an ἄβαξ, a table covered in dust or sand (or perhaps very small rocks). Over the next several centuries, the Greek *abax* evolved into the medieval European *abacus*.

### Multiplication by duplation and mediation

Here is an even older algorithm for multiplying large numbers, sometimes called (*Russian*) *peasant multiplication*. A variant of this method was copied into the Rhind papyrus by the Egyptian scribe Ahmes around 1650 BC, from a document he claimed was (then) about 350 years old. This was the most common method of calculation by Europeans before Fibonacci's introduction of Arabic numerals; it was still taught in elementary schools in Eastern Europe in the late 20th century. This algorithm was also commonly used by early digital computers that did not implement integer multiplication directly in hardware.

<code>PEASANTMULTIPLY(<math>x, y</math>):</code>	$x$	$y$	$prod$
$prod \leftarrow 0$			0
while $x > 0$			
if $x$ is odd	123	+456	= 456
$prod \leftarrow prod + y$	61	+912	= 1368
$x \leftarrow \lfloor x/2 \rfloor$	30	1824	
$y \leftarrow y + y$	15	+3648	= 5016
return $p$	7	+7296	= 12312
	3	+14592	= 26904
	1	+29184	= 56088

The peasant multiplication algorithm breaks the difficult task of general multiplication into four simpler operations: (1) determining parity (even or odd), (2) addition, (3) duplation (doubling a number), and (4) mediation (halving a number, rounding down).<sup>5</sup> Of course a full specification of this algorithm requires describing how to perform those four ‘primitive’ operations. Peasant multiplication requires (a constant factor!) more paperwork to execute by hand, but the necessary operations are easier (for humans) to remember than the  $10 \times 10$  multiplication table required by the American grade school algorithm.<sup>6</sup>

The correctness of peasant multiplication follows from the following recursive identity, which holds for any non-negative integers  $x$  and  $y$ :

$$x \cdot y = \begin{cases} 0 & \text{if } x = 0 \\ \lfloor x/2 \rfloor \cdot (y + y) & \text{if } x \text{ is even} \\ \lfloor x/2 \rfloor \cdot (y + y) + y & \text{if } x \text{ is odd} \end{cases}$$

### Congressional Apportionment

Here is another good example of an algorithm that comes from outside the world of computing. Article I, Section 2 of the United States Constitution requires that

Representatives and direct Taxes shall be apportioned among the several States which may be included within this Union, according to their respective Numbers. . . . The Number of Representatives shall not exceed one for every thirty Thousand, but each State shall have at Least one Representative. . . .

Since there are a limited number of seats available in the House of Representatives, exact proportional representation is impossible without either shared or fractional representatives, neither of which are

<sup>5</sup>The version of this algorithm actually used in ancient Egypt does not use mediation or parity, but it does use comparisons. To avoid halving, the algorithm pre-computes two tables by repeated doubling: one containing all the powers of 2 not exceeding  $x$ , the other containing the same powers of 2 multiplied by  $y$ . The powers of 2 that sum to  $x$  are then found by greedy subtraction, and the corresponding entries in the other table are added together to form the product.

<sup>6</sup>American school kids learn a variant of the *lattice* multiplication algorithm developed by Indian mathematicians and described by Fibonacci in *Liber Abaci*. The two algorithms are equivalent if the input numbers are represented in binary.

legal. As a result, several different apportionment algorithms have been proposed and used to round the fractional solution fairly. The algorithm actually used today, called *the Huntington-Hill method* or *the method of equal proportions*, was first suggested by Census Bureau statistician Joseph Hill in 1911, refined by Harvard mathematician Edward Huntington in 1920, adopted into Federal law ([2 U.S.C. §§2a and 2b](#)) in 1941, and survived a Supreme Court challenge in 1992.<sup>7</sup> The input array  $\text{Pop}[1..n]$  stores the populations of the  $n$  states, and  $R$  is the total number of representatives. Currently,  $n = 50$  and  $R = 435$ . The output array  $\text{Rep}[1..n]$  stores the number of representatives assigned to each state.

```
APPORTIONCONGRESS( $\text{Pop}[1..n], R$ ):
   $PQ \leftarrow \text{NEWPRIORITYQUEUE}$ 
  for  $i \leftarrow 1$  to  $n$ 
     $\text{Rep}[i] \leftarrow 1$ 
     $\text{INSERT}(PQ, i, \text{Pop}[i]/\sqrt{2})$ 
     $R \leftarrow R - 1$ 
  while  $R > 0$ 
     $s \leftarrow \text{EXTRACTMAX}(PQ)$ 
     $\text{Rep}[s] \leftarrow \text{Rep}[s] + 1$ 
     $\text{INSERT}\left(PQ, s, \text{Pop}[s] / \sqrt{\text{Rep}[s](\text{Rep}[s] + 1)}\right)$ 
     $R \leftarrow R - 1$ 
  return  $\text{Rep}[1..n]$ 
```

This pseudocode description assumes that you know how to implement a priority queue that supports the operations `NEWPRIORITYQUEUE`, `INSERT`, and `EXTRACTMAX`. (The actual law doesn't assume that, of course.) The output of the algorithm, and therefore its correctness, does not depend *at all* on how the priority queue is implemented. The Census Bureau uses an unsorted array, stored in a column of an Excel spreadsheet; you should have learned a more efficient solution in your undergraduate data structures class.

## A bad example

Consider “Martin’s algorithm”:<sup>8</sup>

```
BECOMEAMILLIONAIREANDNEVERPAYTAXES:
  Get a million dollars.
  Don't pay taxes.
  If you get caught,
    Say "I forgot."
```

Pretty simple, except for that first step; it's a doozy. A group of billionaire CEOs might consider this an algorithm, since for them the first step is both unambiguous and trivial, but for the rest of us poor slobs, Martin's procedure is too vague to be considered an actual algorithm. On the other hand, this is a perfect example of a **reduction**—it *reduces* the problem of being a millionaire and never paying taxes to the ‘easier’ problem of acquiring a million dollars. We'll see reductions over and over again in this class.

<sup>7</sup>Overruling an earlier ruling by a federal district court, the Supreme Court unanimously held that *any* apportionment method adopted in good faith by Congress is constitutional (*United States Department of Commerce v. Montana*). The current congressional apportionment algorithm is described in gruesome detail at the U.S. Census Department web site <http://www.census.gov/population/www/censusdata/apportionment/computing.html>. A good history of the apportionment problem can be found at <http://www.thirty-thousand.org/pages/Apportionment.htm>. A report by the Congressional Research Service describing various apportionment methods is available at <http://www.rules.house.gov/archives/RL31074.pdf>.

<sup>8</sup>S. Martin, “You Can Be A Millionaire”, Saturday Night Live, January 21, 1978. Appears on *Comedy Is Not Pretty*, Warner Bros. Records, 1979.

As hundreds of businessmen and politicians have demonstrated, if you know how to solve the easier problem, a reduction tells you how to solve the harder one.

Martin's algorithm, like many of our previous examples, is not the kind of algorithm that computer scientists are used to thinking about, because it is phrased in terms of operations that are difficult for computers to perform. In this class, we'll focus (almost!) exclusively on algorithms that can be reasonably implemented on a computer. In other words, each step in the algorithm must be something that either is directly supported by common programming languages (such as arithmetic, assignments, loops, or recursion) or is something that you've already learned how to do in an earlier class (like sorting, binary search, or depth first search).

### 0.3 Writing down algorithms

Computer programs are concrete representations of algorithms, but algorithms are *not* programs; they should not be described in a particular programming language. The whole *point* of this course is to develop computational techniques that can be used in *any* programming language. The idiosyncratic syntactic details of C, C++, C#, Java, Python, Ruby, Erlang, Haskell, OcaML, Scheme, Visual Basic, Smalltalk, Javascript, Processing, Squeak, Forth, TeX, Fortran, COBOL, Intercal, or Brainfuck are of little or no importance in algorithm design, and focusing on them will only distract you from what's *really* going on.<sup>9</sup> What we really want is closer to what you'd write in the *comments* of a real program than the code itself.

On the other hand, a plain English prose description is usually not a good idea either. Algorithms have a lot of structure—especially conditionals, loops, and recursion—that are far too easily hidden by unstructured prose. Like any natural language, English is full of ambiguities, subtleties, and shades of meaning, but algorithms must be described as accurately as possible. Finally and more seriously, there is natural tendency to describe repeated operations informally: “Do this first, then do this second, and so on.” But as anyone who has taken one of those ‘What comes next in this sequence?’ tests already knows, specifying what happens in the first few iterations of a loop says very little about what happens later iterations. To make the description unambiguous, we must explicitly specify the behavior of *every* iteration.

The best way to write down an algorithm is using pseudocode. Pseudocode uses the *structure* of formal programming languages and mathematics to break algorithms into primitive steps; but the primitive steps themselves may be written using mathematics, pure English, or an appropriate mixture of the two. Well-written pseudocode reveals the internal structure of the algorithm but hides irrelevant implementation details, making the algorithm much easier to understand, analyze, debug, and implement.

The precise syntax of pseudocode is a personal choice, but the overriding goal should be clarity and precision. Ideally, pseudocode should allow any competent programmer to implement the underlying algorithm, quickly and correctly, in *their* favorite programming language, *without understanding why the algorithm works*. Here are the guidelines I follow and strongly recommend:

---

<sup>9</sup>This is, of course, a matter of religious conviction. Linguists argue incessantly over the *Sapir-Whorf hypothesis*, which states (more or less) that people think only in the categories imposed by their languages. According to an extreme formulation of this principle, some concepts in one language simply cannot be understood by speakers of other languages, not just because of technological advancement—How would you translate ‘jump the shark’ or ‘blog’ into Aramaic?—but because of inherent structural differences between languages and cultures. For a more skeptical view, see Steven Pinker’s *The Language Instinct*. There is admittedly some strength to this idea when applied to different programming paradigms. (What’s the Y combinator, again? How do templates work? What’s an Abstract Factory?) Fortunately, those differences are generally too subtle to have much impact in *this* class.

- Be consistent!
- Use standard imperative programming keywords (if/then/else, while, for, repeat/until, case, return) and notation (*variable*  $\leftarrow$  *value*, *Array*[*index*], *function(argument)*, *bigger* > *smaller*, etc.). Keywords should be standard English words: write ‘else if’ instead of ‘elif’.
- Indent everything carefully and consistently; the block structure should be visible from across the room. This rule is especially important for nested loops and conditionals. *Don’t* add unnecessary syntactic sugar like braces or begin/end tags; careful indentation is almost always enough.
- Use mnemonic algorithm and variable names. Short variable names are good, but readability is more important than concision; except for idioms like loop indices, short but complete words are better than single letters. Absolutely *never* use pronouns!
- Use standard mathematical notation for standard mathematical things. For example, write  $x \cdot y$  instead of  $x * y$  for multiplication; write  $x \bmod y$  instead of  $x \% y$  for remainder; write  $\sqrt{x}$  instead of  $\text{sqrt}(x)$  for square roots; write  $a^b$  instead of  $\text{power}(a, b)$  for exponentiation; and write  $\phi$  instead of  $\text{phi}$  for the golden ratio.
- Avoid mathematical notation if English is clearer. For example, ‘Insert  $a$  into  $X$ ’ may be preferable to  $\text{INSERT}(X, a)$  or  $X \leftarrow X \cup \{a\}$ .
- Each statement should fit on one line, and each line should contain either exactly one statement or exactly one structuring element (for, while, if). (I sometimes make an exception for short and similar statements like  $i \leftarrow i + 1$ ;  $j \leftarrow j - 1$ ;  $k \leftarrow 0$ .)
- *Don’t* use a fixed-width typeface to typeset pseudocode; it’s much harder to read than normal typeset text. Similarly, *don’t* typeset keywords like ‘for’ or ‘while’ in a different **style**; the syntactic sugar is not what you want the reader to look at. On the other hand, I use *italics* for variables, **Small Caps** for algorithms and constants, and *a different typeface* for literal strings and comments.

## 0.4 Analyzing algorithms

It’s not enough just to write down an algorithm and say ‘Behold!’ We also need to convince ourselves (and our graders) that the algorithm does what it’s supposed to do, and that it does it efficiently.

### Correctness

In some application settings, it is acceptable for programs to behave correctly most of the time, on all ‘reasonable’ inputs. Not in this class; we require algorithms that are correct for *all possible* inputs. Moreover, we must *prove* that our algorithms are correct; trusting our instincts, or trying a few test cases, isn’t good enough. Sometimes correctness is fairly obvious, especially for algorithms you’ve seen in earlier courses. On the other hand, ‘obvious’ is all too often a synonym for ‘wrong’. Many of the algorithms we will discuss in this course will require extra work to prove correct. Correctness proofs almost always involve induction. We *like* induction. Induction is our friend.<sup>10</sup>

But before we can formally prove that our algorithm does what we want it to do, we have to formally *state* what we want it to do! Usually problems are given to us in real-world terms, not with formal mathematical descriptions. It’s up to us, the algorithm designers, to restate these problems in terms of mathematical objects that we can prove things about—numbers, arrays, lists, graphs, trees, and so on.

---

<sup>10</sup>If induction is *not* your friend, you will have a hard time in this course.

We also need to determine if the problem statement makes any hidden assumptions, and state those assumptions explicitly. (For example, in the song “*n Bottles of Beer on the Wall*”,  $n$  is always a positive integer.) Restating the problem formally is not only required for proofs; it is also one of the best ways to really understand what the problems is asking for. The hardest part of answering any question is figuring out the right way to ask it!

It is important to remember the distinction between a problem and an algorithm. A problem is a task to perform, like “Compute the square root of  $x$ ” or “Sort these  $n$  numbers” or “Keep  $n$  algorithms students awake for  $t$  minutes”. An algorithm is a set of instructions for accomplishing such a task. The same problem may have hundreds of different algorithms; the same algorithm may solve hundreds of different problems.

### Running time

The most common way of ranking different algorithms for the same problem is by how fast they run. Ideally, we want the fastest possible algorithm for our problem. In many application settings, it is acceptable for programs to run efficiently most of the time, on all ‘reasonable’ inputs. Not in this class; we require algorithms that *always* run efficiently, even in the worst case.

But how do we measure running time? As a specific example, how long does it take to sing the song *BOTTLESOFBEER( $n$ )*? This is obviously a function of the input value  $n$ , but it also depends on how quickly you can sing. Some singers might take ten seconds to sing a verse; others might take twenty. Technology widens the possibilities even further. Dictating the song over a telegraph using Morse code might take a full minute per verse. Downloading an mp3 over the Web might take a tenth of a second per verse. Duplicating the mp3 in a computer’s main memory might take only a few microseconds per verse.

What’s important here is how the singing time changes as  $n$  grows. Singing *BOTTLESOFBEER( $2n$ )* takes about twice as long as singing *BOTTLESOFBEER( $n$ )*, no matter what technology is being used. This is reflected in the asymptotic singing time  $\Theta(n)$ . We can measure time by counting how many times the algorithm executes a certain instruction or reaches a certain milestone in the ‘code’. For example, we might notice that the word ‘beer’ is sung three times in every verse of *BOTTLESOFBEER*, so the number of times you sing ‘beer’ is a good indication of the total singing time. For this question, we can give an exact answer: *BOTTLESOFBEER( $n$ )* uses exactly  $3n + 3$  beers.

There are plenty of other songs that have non-trivial singing time. This one is probably familiar to most English-speakers:

```
NDAYSOFCHRISTMAS(gifts[2..n]):
for i ← 1 to n
    Sing “On the ith day of Christmas, my true love gave to me”
    for j ← i down to 2
        Sing “j gifts[j]”
    if i > 1
        Sing “and”
    Sing “a partridge in a pear tree.”
```

The input to *NDAYSOFCHRISTMAS* is a list of  $n - 1$  gifts. It’s quite easy to show that the singing time is  $\Theta(n^2)$ ; in particular, the singer mentions the name of a gift  $\sum_{i=1}^n i = n(n + 1)/2$  times (counting the partridge in the pear tree). It’s also easy to see that during the first  $n$  days of Christmas, my true love gave to me exactly  $\sum_{i=1}^n \sum_{j=1}^i j = n(n + 1)(n + 2)/6 = \Theta(n^3)$  gifts. Other songs that take quadratic time to sing are “Old MacDonald Had a Farm”, “There Was an Old Lady Who Swallowed a Fly”, “The House that Jack Built”, “Hole in the Bottom of the Sea”, “Green Grow the Rushes O”, “Eh, Cumpari!”, “Alouette”, “Echad Mi Yode’a”, “Chad Gadya”, “Minkurinn í hænsnakofanum”, and “Ist das nicht ein Schnitzelbank?” For further details, consult your nearest preschooler.

```
OLDMACDONALD(animals[1 .. n],noise[1 .. n]):
  for i  $\leftarrow$  1 to n
    Sing "Old MacDonald had a farm, E I E I O"
    Sing "And on this farm he had some animals[i], E I E I O"
    Sing "With a noise[i] noise[i] here, and a noise[i] noise[i] there"
    Sing "Here a noise[i], there a noise[i], everywhere a noise[i] noise[i]"
    for j  $\leftarrow$  i - 1 down to 1
      Sing "noise[j] noise[j] here, noise[j] noise[j] there"
      Sing "Here a noise[j], there a noise[j], everywhere a noise[j] noise[j]"
    Sing "Old MacDonald had a farm, E I E I O."
```

```
ALOUETTE(lapart[1 .. n]):
  Chantez « Alouette, gentille alouette, alouette, je te plumerais. »
  pour tout i de 1 à n
    Chantez « Je te plumerais lapart[i]. Je te plumerais lapart[i]. »
    pour tout j de i - 1 à bas à 1
      Chantez « Et lapart[j] ! Et lapart[j] ! »
    Chantez « Ooooooo ! »
  Chantez « Alouette, gentille alluette, alouette, je te plumerais. »
```

For a slightly more complicated example, consider the algorithm APPORTIONCONGRESS. Here the running time obviously depends on the implementation of the priority queue operations, but we can certainly bound the running time as  $O(N + RI + (R - n)E)$ , where  $N$  denotes the running time of NEWPRIORITYQUEUE,  $I$  denotes the running time of INSERT, and  $E$  denotes the running time of EXTRACTMAX. Under the reasonable assumption that  $R > 2n$  (on average, each state gets at least two representatives), we can simplify the bound to  $O(N + R(I + E))$ . The Census Bureau implements the priority queue using an unsorted array of size  $n$ ; this implementation gives us  $N = I = \Theta(1)$  and  $E = \Theta(n)$ , so the overall running time is  $O(Rn)$ . This is good enough for government work, but we can do better. Implementing the priority queue using a binary heap (or a heap-ordered array) gives us  $N = \Theta(1)$  and  $I = R = O(\log n)$ , which implies an overall running time of  $O(R \log n)$ .

Sometimes we are also interested in other computational resources: space, randomness, page faults, inter-process messages, and so forth. We can use the same techniques to analyze those resources as we use to analyze running time.

## 0.5 A Longer Example: Stable Matching

Every year, thousands of new doctors must obtain internships at hospitals around the United States. During the first half of the 20th century, competition among hospitals for the best doctors led to earlier and earlier offers of internships, sometimes as early as the second year of medical school, along with tighter deadlines for acceptance. In the 1940s, medical schools agreed not to release information until a common date during their students' fourth year. In response, hospitals began demanding faster decisions. By 1950, hospitals would regularly call doctors, offer them internships, and demand *immediate* responses. Interns were forced to gamble if their third-choice hospital called first—accept and risk losing a better opportunity later, or reject and risk having no position at all.<sup>11</sup>

Finally, a central clearinghouse for internship assignments, now called the National Resident Matching Program, was established in the early 1950s. Each year, doctors submit a ranked list of all hospitals where they would accept an internship, and each hospital submits a ranked list of doctors they would

<sup>11</sup>The academic job market involves similar gambles, at least in computer science. Some departments start making offers in February with two-week decision deadlines; other departments don't even start interviewing until late March; MIT notoriously waits until May, when all its interviews are over, before making *any* faculty offers.

accept as interns. The NRMP then computes an assignment of interns to hospitals that satisfies the following *stability* requirement. For simplicity, let's assume that there are  $n$  doctors and  $n$  hospitals; each hospital offers exactly one internship; each doctor ranks all hospitals and vice versa; and finally, there are no ties in the doctors' and hospitals' rankings.<sup>12</sup> We say that a matching of doctors to hospitals is *unstable* if there are two doctors  $\alpha$  and  $\beta$  and two hospitals  $A$  and  $B$ , such that

- $\alpha$  is assigned to  $A$ , and  $\beta$  is assigned to  $B$ ;
- $\alpha$  prefers  $B$  to  $A$ , and  $B$  prefers  $\alpha$  to  $\beta$ .

In other words,  $\alpha$  and  $B$  would both be happier with each other than with their current assignment. The goal of the Resident Match is a *stable matching*, in which no doctor or hospital has an incentive to cheat the system. At first glance, it is not clear that a stable matching exists!

In 1952, the NRMP adopted the “Boston Pool” algorithm to assign interns, so named because it had been previously used by a regional clearinghouse in the Boston area. The algorithm is often inappropriately attributed to David Gale and Lloyd Shapley, who formally analyzed the algorithm and first proved that it computes a stable matching in 1962; Gale and Shapley used the metaphor of college admissions.<sup>13</sup> Similar algorithms have since been adopted for other matching markets, including faculty recruiting in France, university admission in Germany, public school admission in New York and Boston, and billet assignments for US Navy sailors. The stable matching problem is

The Boston Pool algorithm proceeds in rounds until every position has been filled. Each round has two stages:

1. An arbitrary unassigned hospital  $A$  offers its position to the best doctor  $\alpha$  (according to the hospital's preference list) who has not already rejected it.
2. Each doctor plans to accept the best hospital (according to her preference list) that makes her an offer. Thus, if  $\alpha$  is currently unassigned, she (tentatively) accepts the offer from  $A$ . If  $\alpha$  already has an assignment but prefers  $A$ , she rejects her existing assignment and (tentatively) accepts the new offer from  $A$ . Otherwise,  $\alpha$  rejects the new offer.

For example, suppose four doctors (Dr. Quincy, Dr. Rotwang, Dr. Shephard, and Dr. Tam, represented by lower-case letters) and four hospitals (Arkham Asylum, Bethlehem Royal Hospital, County General Hospital, and The Dharma Initiative, represented by upper-case letters) rank each other as follows:

$q$	$r$	$s$	$t$	$A$	$B$	$C$	$D$
$A$	A	B	D	t	r	t	s
$B$	D	A	B	s	t	q	r
$C$	C	C	C	r	q	r	q
$D$	B	D	A	q	s	s	t

Given these preferences as input, the Boston Pool algorithm might proceed as follows:

<sup>12</sup>In reality, most hospitals offer multiple internships, *each doctor ranks only a subset of the hospitals and vice versa*, and there are typically more internships than interested doctors. And then it starts getting complicated.

<sup>13</sup>The “Gale-Shapley algorithm” is a prime instance of *Stigler's Law of Eponymy*: *No scientific discovery is named after its original discoverer*. In his 1980 paper that gives the law its name, the statistician Stephen Stigler claimed that this law was first proposed by sociologist Robert K. Merton. However, similar statements were previously made by Vladimir Arnold in the 1970's (“Discoveries are rarely attributed to the correct person.”), Carl Boyer in 1968 (“Clio, the muse of history, often is fickle in attaching names to theorems!”), Alfred North Whitehead in 1917 (“Everything of importance has been said before by someone who did not discover it.”), and even Stephen's father George Stigler in 1966 (“If we should ever encounter a case where a theory is named for the correct man, it will be noted.”). We will see many other examples of Stigler's law in this class.

1. Arkham makes an offer to Dr. Tam.
2. Bedlam makes an offer to Dr. Rotwang.
3. County makes an offer to Dr. Tam, who rejects her earlier offer from Arkham.
4. Dharma makes an offer to Dr. Shephard. (From this point on, because there is only one unmatched hospital, the algorithm has no more choices.)
5. Arkham makes an offer to Dr. Shephard, who rejects her earlier offer from Dharma.
6. Dharma makes an offer to Dr. Rotwang, who rejects her earlier offer from Bedlam.
7. Bedlam makes an offer to Dr. Tam, who rejects her earlier offer from County.
8. County makes an offer to Dr. Quincy.

At this point, all pending offers are accepted, and the algorithm terminates with a matching:  $(A, s)$ ,  $(B, t)$ ,  $(C, q)$ ,  $(D, r)$ . You can (and should) verify by brute force that this matching is stable, even though no doctor was hired by her favorite hospital, and no hospital hired its favorite doctor. In fact, this is **the only stable matching** for this list of preferences.

Analyzing the algorithm's running time is relatively straightforward. Each hospital makes an offer to each doctor at most once, so the algorithm requires at most  $n^2$  rounds. In an actual implementation, each doctor and hospital can be identified by a unique integer between 1 and  $n$ , and the preference lists can be represented as two arrays  $DocPref[1..n][1..n]$  and  $HosPref[1..n][1..n]$ , where  $DocPref[\alpha][r]$  represents the  $r$ th hospital in doctor  $\alpha$ 's preference list, and  $HosPref[A][r]$  represents the  $r$ th doctor in hospital  $A$ 's preference list. With the input in this form, the Boston Pool algorithm can be implemented to run in  **$O(n^2)$  time**; we leave the details as **an easy exercise**.

A somewhat harder exercise is to prove that there are inputs (and choices of who makes offers when) that force  $\Omega(n^2)$  rounds before the algorithm terminates. Thus, the  $O(n^2)$  upper bound on the worst-case running time cannot be improved; in this case, we say our analysis is **tight**.

## Correctness

But why is the algorithm *correct*? Gale and Shapley proved that the Boston Pool algorithm always computes a stable matching as follows. The algorithm continues as long as there is at least one unfilled position; conversely, when the algorithm terminates (after at most  $n^2$  rounds), every position is filled. No doctor can accept more than one position, and no hospital can hire more than one doctor. Thus, the algorithm always computes a matching; it remains only to prove that the matching is stable.

Suppose doctor  $\alpha$  is assigned to hospital  $A$  in the final matching, but prefers  $B$ . Because every doctor accepts the best offer she receives,  $\alpha$  received no offer she liked more than  $A$ . In particular,  $B$  never made an offer to  $\alpha$ . On the other hand,  $B$  made offers to every doctor they like more than  $\beta$ . Thus,  $B$  prefers  $\beta$  to  $\alpha$ , and so there is no instability.

Surprisingly, the correctness of the algorithm does not depend on which hospital makes its offer in which round. In fact, there is a stronger sense in which the order of offers doesn't matter—no matter which unassigned hospital makes an offer in each round, *the algorithm always computes the same matching!* Let's say that  $\alpha$  is a **feasible** doctor for  $A$  if there is a stable matching that assigns doctor  $\alpha$  to hospital  $A$ .

**Lemma 1.** *During the Boston Pool algorithm, each hospital  $A$  is rejected only by doctors that are infeasible for  $A$ .*

**Proof:** We prove the lemma by induction. Consider an arbitrary round of the Boston Pool algorithm, in which doctor  $\alpha$  rejects one hospital  $A$  for another hospital  $B$ . The rejection implies that  $\alpha$  prefers  $B$  to  $A$ .

Every doctor that appears higher than  $\alpha$  in  $B$ 's preference list has already rejected  $B$  and therefore, by the inductive hypothesis, is infeasible for  $B$ .

Now consider an arbitrary matching that assigns  $\alpha$  to  $A$ . We already established that  $\alpha$  prefers  $B$  to  $A$ . If  $B$  prefers  $\alpha$  to its partner, the matching is unstable. On the other hand, if  $B$  prefers its partner to  $\alpha$ , then (by our earlier argument) its partner is infeasible, and again the matching is unstable. We conclude that there is no stable matching that assigns  $\alpha$  to  $A$ .  $\square$

Now let  $\text{best}(A)$  denote the highest-ranked *feasible* doctor on  $A$ 's preference list. Lemma 1 implies that every doctor that  $A$  prefers to its final assignment is infeasible for  $A$ . On the other hand, the final matching is stable, so the doctor assigned to  $A$  is feasible for  $A$ . The following result is now immediate:

**Corollary 1.1.** *The Boston Pool algorithm assigns  $\text{best}(A)$  to  $A$ , for every hospital  $A$ .*

Thus, from the hospitals' point of view, the Boston Pool algorithm computes the best possible stable matching. It turns out that this is also the *worst* possible matching from the doctors' viewpoint! Let  $\text{worst}(\alpha)$  denote the lowest-ranked feasible hospital on doctor  $\alpha$ 's preference list.

**Corollary 1.2.** *The Boston Pool algorithm assigns  $\alpha$  to  $\text{worst}(\alpha)$ , for every doctor  $\alpha$ .*

**Proof:** Suppose the Boston Pool algorithm assigns doctor  $\alpha$  to hospital  $A$ ; we need to show that  $A = \text{worst}(\alpha)$ .

Consider an arbitrary stable matching where  $A$  is *not* matched with  $\alpha$  but with another doctor  $\beta$ . The previous corollary implies that  $A$  prefers  $\alpha = \text{best}(A)$  to  $\beta$ . Because the matching is stable,  $\alpha$  must therefore prefer her assigned hospital to  $A$ . This argument works for *any* stable assignment, so  $\alpha$  prefers *every* other feasible match to  $A$ ; in other words,  $A = \text{worst}(\alpha)$ .  $\square$

A subtle consequence of these two corollaries, discovered by Dubins and Freeman in 1981, is that a doctor can potentially improve her assignment by lying about her preferences, but a hospital cannot. (However, a set of hospitals can collude so that *some* of their assignments improve.) Partly for this reason, the National Residency Matching Program reversed its matching algorithm in 1998, so that potential residents offer to work for hospitals in preference order, and each hospital accepts its best offer. Thus, the new algorithm computes the best possible stable matching for the doctors, and the worst possible stable matching for the hospitals. In practice, however, this modification affected less than 1% of the resident's assignments. As far as I know, the precise effect of this change on the *patients* is an open problem.

## 0.6 Why are we here, anyway?

This class is ultimately about learning two skills that are crucial for all computer scientists.

1. **Intuition:** How to *think* about abstract computation.
2. **Language:** How to *talk* about abstract computation.

The first goal of this course is to help you develop algorithmic *intuition*. How do various algorithms really work? When you see a problem for the first time, how should you attack it? How do you tell which techniques will work at all, and which ones will work best? How do you judge whether one algorithm is better than another? How do you tell whether you have the best possible solution? These are *not* easy questions; anyone who says differently is selling something.

Our second main goal is to help you develop algorithmic *language*. It's not enough just to understand how to solve a problem; you also have to be able to explain your solution to somebody else. I don't mean just how to turn your algorithms into working code—despite what many students (and inexperienced programmers) think, ‘somebody else’ is *not* just a computer. Nobody programs alone. Code is read far more often than it is written, or even compiled. Perhaps more importantly in the short term, explaining something to somebody else is one of the best ways to clarify your own understanding. As Albert Einstein (or was it Richard Feynman?) apocryphally put it, “You do not really understand something unless you can explain it to your grandmother.”

Along the way, you'll pick up a bunch of algorithmic facts—mergesort runs in  $\Theta(n \log n)$  time; the amortized time to search in a splay tree is  $O(\log n)$ ; greedy algorithms usually don't produce optimal solutions; the traveling salesman problem is NP-hard—but these aren't the point of the course. You can always look up mere facts in a textbook or on the web, provided you have enough intuition and experience to know what to look for. That's why we let you bring cheat sheets to the exams; we don't want you wasting your study time trying to memorize all the facts you've seen.

You'll also practice a lot of algorithm design and analysis skills—finding useful (counter)examples, developing induction proofs, solving recurrences, using big-Oh notation, using probability, giving problems crisp mathematical descriptions, and so on. These skills are *incredibly* useful, and it's impossible to develop good intuition and good communication skills without them, but they aren't the main point of the course either. At this point in your educational career, you should be able to pick up most of those skills on your own, once you know what you're trying to do.

Unfortunately, there is no systematic procedure—no algorithm—to determine which algorithmic techniques are most effective at solving a given problem, or finding good ways to explain, analyze, optimize, or implement a given algorithm. Like many other activities (music, writing, juggling, acting, martial arts, sports, cooking, programming, teaching, etc.), the *only* way to master these skills is to make them your own, through practice, practice, and more practice. You can only develop good problem-solving skills by solving problems. You can only develop good communication skills by communicating. Good intuition is the product of experience, not its replacement. We *can't* teach you how to do well in this class. All we can do (and what we will do) is lay out some fundamental tools, show you how to use them, create opportunities for you to practice with them, and give you honest feedback, based on our own hard-won experience and intuition. The rest is up to you.

Good algorithms are extremely useful, elegant, surprising, deep, even beautiful, but most importantly, algorithms are *fun!* I hope you will enjoy playing with them as much as I do.



Boethius the algorist versus Pythagoras the abacist.  
from *Margarita Philosophica* by Gregor Reisch (1503)

## Exercises

0. Describe and analyze an algorithm that determines, given a legal arrangement of standard pieces on a standard chess board, which player will win at chess from the given starting position if both players play perfectly. [*Hint: There is a one-line solution!*]
1. The traditional Devonian/Cornish drinking song “The Barley Mow” has the following pseudolyrics, where  $\text{container}[i]$  is the name of a container that holds  $2^i$  ounces of beer. One version of the song uses the following containers: npperkin, gill pot, half-pint, pint, quart, pottle, gallon, half-anker, anker, firkin, half-barrel, barrel, hogshead, pipe, well, river, and ocean. (Every container in this list is twice as big as its predecessor, except that a firkin is actually 2.25 ankers, and the last three units are just silly.)

```
BARLEYMOW( $n$ ):
  "Here's a health to the barley-mow, my brave boys,"
  "Here's a health to the barley-mow!"

  "We'll drink it out of the jolly brown bowl,"
  "Here's a health to the barley-mow!"
  "Here's a health to the barley-mow, my brave boys,"
  "Here's a health to the barley-mow!"

  for  $i \leftarrow 1$  to  $n$ 
    "We'll drink it out of the container[ $i$ ], boys,"
    "Here's a health to the barley-mow!"
    for  $j \leftarrow i$  downto 1
      "The container[j],"
      "And the jolly brown bowl!"
      "Here's a health to the barley-mow!"
      "Here's a health to the barley-mow, my brave boys,"
      "Here's a health to the barley-mow!"
```

- (a) Suppose each container name  $\text{container}[i]$  is a single word, and you can sing four words a second. How long would it take you to sing  $\text{BARLEYMOW}(n)$ ? (Give a tight asymptotic bound.)
- (b) If you want to sing this song for  $n > 20$ , you'll have to make up your own container names. To avoid repetition, these names will get progressively longer as  $n$  increases<sup>14</sup>. Suppose  $\text{container}[n]$  has  $\Theta(\log n)$  syllables, and you can sing six syllables per second. Now how long would it take you to sing  $\text{BARLEYMOW}(n)$ ? (Give a tight asymptotic bound.)
- (c) Suppose each time you mention the name of a container, you actually drink the corresponding amount of beer: one ounce for the jolly brown bowl, and  $2^i$  ounces for each  $\text{container}[i]$ . Assuming for purposes of this problem that you are at least 21 years old, exactly how many ounces of beer would you drink if you sang  $\text{BARLEYMOW}(n)$ ? (Give an exact answer, not just an asymptotic bound.)
2. Describe and analyze the Boston Pool stable matching algorithm in more detail, so that the worst-case running time is  $O(n^2)$ , as claimed earlier in the notes.

<sup>14</sup>“We'll drink it out of the hemisemidemiyottapint, boys!”

3. Prove that it is possible for the Boston Pool algorithm to execute  $\Omega(n^2)$  rounds. (You need to describe both a suitable input and a sequence of  $\Omega(n^2)$  valid proposals.)
4. Describe and analyze an efficient algorithm to determine whether a given set of hospital and doctor preferences has to a *unique* stable matching.
5. Consider a generalization of the stable matching problem, where some doctors do not rank all hospitals and some hospitals do not rank all doctors, and a doctor can be assigned to a hospital only if each appears in the other's preference list. In this case, there are three additional unstable situations:
  - A hospital prefers an unmatched doctor to its assigned match.
  - A doctor prefers an unmatched hospital to her assigned match.
  - An unmatched doctor and an unmatched hospital appear in each other's preference lists.

Describe and analyze an efficient algorithm that computes a stable matching in this setting.

Note that a stable matching may leave some doctors and hospitals unmatched, even though their preference lists are non-empty. For example, if every doctor lists Harvard as their only acceptable hospital, and every hospital lists Dr. House as their only acceptable intern, then only House and Harvard will be matched.

6. Recall that the input to the Huntington-Hill apportionment algorithm APPORTIONCONGRESS is an array  $P[1..n]$ , where  $P[i]$  is the population of the  $i$ th state, and an integer  $R$ , the total number of representatives to be allotted. The output is an array  $r[1..n]$ , where  $r[i]$  is the number of representatives allotted to the  $i$ th state by the algorithm.

Let  $P = \sum_{i=1}^n P[i]$  denote the total population of the country, and let  $r_i^* = R \cdot P[i]/P$  denote the ideal number of representatives for the  $i$ th state.

- (a) Prove that  $r[i] \geq \lfloor r_i^* \rfloor$  for all  $i$ .
- (b) Describe and analyze an algorithm that computes exactly the same congressional apportionment as APPORTIONCONGRESS in  $O(n \log n)$  time. (Recall that the running time of APPORTIONCONGRESS depends on  $R$ , which could be arbitrarily larger than  $n$ .)
- (c) If a state's population is small relative to the other states, its ideal number  $r_i^*$  of representatives could be close to zero; thus, tiny states are over-represented by the Huntington-Hill apportionment process. Surprisingly, this can also be true of very large states. Let  $\alpha = (1 + \sqrt{2})/2 \approx 1.20710678119$ . Prove that for any  $\varepsilon > 0$ , there is an input to APPOINTCONGRESS with  $\max_i P[i] = P[1]$ , such that  $r[1] > (\alpha - \varepsilon)r_1^*$ .
- ★(d) Can you improve the constant  $\alpha$  in the previous question?

*Our life is frittered away by detail. Simplify, simplify.*

— Henry David Thoreau

*The control of a large force is the same principle as the control of a few men:  
it is merely a question of dividing up their numbers.*

— Sun Zi, *The Art of War* (c. 400 C.E.), translated by Lionel Giles (1910)

*Nothing is particularly hard if you divide it into small jobs.*

— Henry Ford

## 1 Recursion

### 1.1 Simplify and delegate

*Reduction* is the single most common technique used in designing algorithms. Reducing one problem  $X$  to another problem (or set of problems)  $Y$  means to write an algorithm for  $X$ , using an algorithm or  $Y$  as a subroutine or black box.

For example, the congressional apportionment algorithm described in Lecture 0 reduces the problem of apportioning Congress to the problem of maintaining a priority queue under the operations `INSERT` and `EXTRACTMAX`. Those data structure operations are black boxes; the apportionment algorithm does not depend on any specific implementation. Conversely, when we design a particular priority queue data structure, we typically neither know nor care how our data structure will be used. Whether or not the Census Bureau plans to use our code to apportion Congress is completely irrelevant to our design choices. As a general rule, when we design algorithms, we may not know—and we should not care—how the basic building blocks we use are implemented, or how your algorithm might be used as a basic building block to solve a bigger problem.

A particularly powerful kind of reduction is *recursion*, which can be defined loosely as follows:

- If the given instance of the problem is small or simple enough, just solve it.
- Otherwise, reduce the problem to one or more *simpler instances of the same problem*.

If the self-reference is confusing, it's helpful to imagine that someone else is going to solve the simpler problems, just as you would assume for other types of reductions. Your *only* task is to *simplify* the original problem, or to solve it directly when simplification is either unnecessary or impossible. The Recursion Fairy will magically take care of the simpler subproblems.<sup>1</sup>

There is one mild technical condition that must be satisfied in order for any recursive method to work correctly, namely, that there is no infinite sequence of reductions to ‘simpler’ and ‘simpler’ subproblems. Eventually, the recursive reductions must stop with an elementary *base case* that can be solved by some other method; otherwise, the recursive algorithm will never terminate. This finiteness condition is almost always satisfied trivially, but we should always be wary of ‘obvious’ recursive algorithms that actually recurse forever.<sup>2</sup>

<sup>1</sup>I used to refer to ‘elves’ instead of the Recursion Fairy, referring to the traditional fairy tale about an old shoemaker who leaves his work unfinished when he goes to bed, only to discover upon waking that elves have finished everything overnight. Someone more entheogenically experienced than I might recognize them as Terence McKenna’s ‘self-adjusting machine elves’.

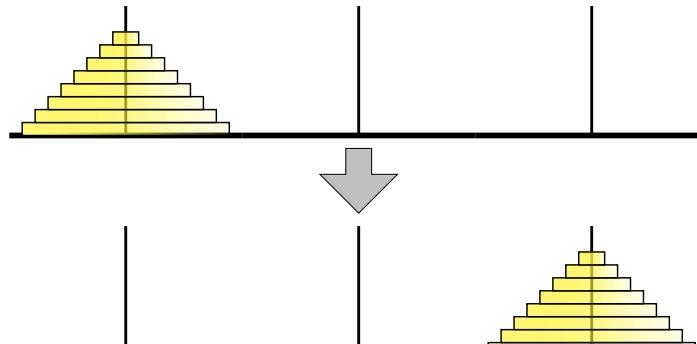
<sup>2</sup>All too often, ‘obvious’ is a synonym for ‘false’.

## 1.2 Tower of Hanoi

The Tower of Hanoi puzzle was first published by the mathematician François Éduard Anatole Lucas in 1883, under the pseudonym ‘N. Claus (de Siam)’ (an anagram of ‘Lucas d’Amiens’). The following year, Henri de Parville described the puzzle with the following remarkable story:<sup>3</sup>

In the great temple at Benares beneath the dome which marks the centre of the world, rests a brass plate in which are fixed three diamond needles, each a cubit high and as thick as the body of a bee. On one of these needles, at the creation, God placed sixty-four discs of pure gold, the largest disc resting on the brass plate, and the others getting smaller and smaller up to the top one. This is the Tower of Bramah. Day and night unceasingly the priests transfer the discs from one diamond needle to another according to the fixed and immutable laws of Bramah, which require that the priest on duty must not move more than one disc at a time and that he must place this disc on a needle so that there is no smaller disc below it. When the sixty-four discs shall have been thus transferred from the needle on which at the creation God placed them to one of the other needles, tower, temple, and Brahmins alike will crumble into dust, and with a thunderclap the world will vanish.

Of course, being good computer scientists, we read this story and immediately substitute  $n$  for the hardwired constant sixty-four. How can we move a tower of  $n$  disks from one needle to another, using a third needles as an occasional placeholder, never placing any disk on top of a smaller disk?



The Tower of Hanoi puzzle

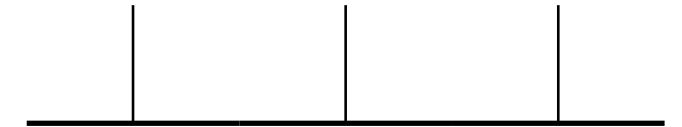
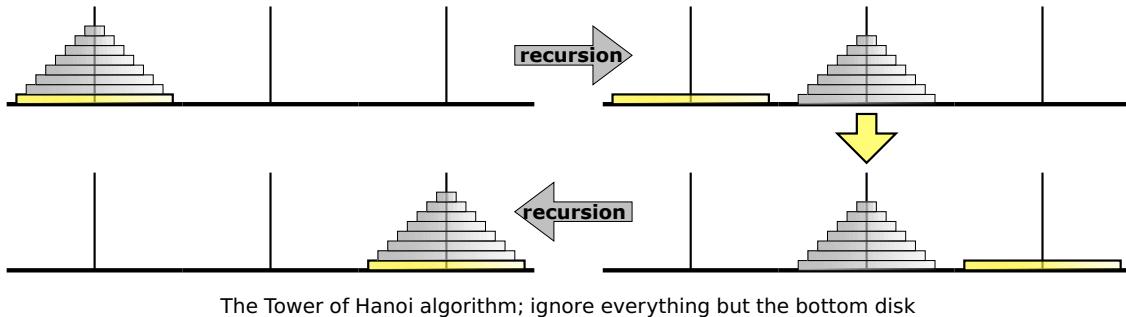
The trick to solving this puzzle is to think recursively. Instead of trying to solve the entire puzzle all at once, let’s concentrate on moving just the largest disk. We can’t move it at the beginning, because all the other disks are covering it; we have to move those  $n - 1$  disks to the third needle before we can move the  $n$ th disk. And then after we move the  $n$ th disk, we have to move those  $n - 1$  disks back on top of it. So now all we have to figure out is how to...

**STOP!!** That’s it! We’re done! We’ve successfully reduced the  $n$ -disk Tower of Hanoi problem to two instances of the  $(n - 1)$ -disk Tower of Hanoi problem, which we can gleefully hand off to the Recursion Fairy (or, to carry the original story further, to the junior monks at the temple).

Our algorithm does make one subtle but important assumption: *There is a largest disk*. In other words, our recursive algorithm works for any  $n \geq 1$ , but it breaks down when  $n = 0$ . We must handle that base case directly. Fortunately, the monks at Benares, being good Buddhists, are quite adept at moving zero disks from one needle to another in no time at all.

While it’s tempting to think about how all those smaller disks get moved—in other words, **what happens when the recursion is unfolded**—it’s not necessary. In fact, for more complicated problems, unfolding the recursive calls is merely distracting. Our *only* task is to reduce the problem to one or more simpler instances, or to solve the problem directly if such a reduction is impossible. Our algorithm is trivially correct when  $n = 0$ . For any  $n \geq 1$ , the Recursion Fairy correctly moves (or more formally, the

<sup>3</sup>This English translation is from W. W. Rouse Ball and H. S. M. Coxeter’s book *Mathematical Recreations and Essays*.



The base case for the Tower of Hanoi algorithm. There is no spoon.

inductive hypothesis implies that our algorithm correctly moves) the top  $n - 1$  disks, so our algorithm is clearly correct.

Here's the recursive Hanoi algorithm in more typical pseudocode.

```
HANOI( $n, src, dst, tmp$ ):
  if  $n > 0$ 
    HANOI( $n - 1, src, tmp, dst$ )
    move disk  $n$  from  $src$  to  $dst$ 
    HANOI( $n - 1, tmp, dst, src$ )
```

Let  $T(n)$  denote the number of moves required to transfer  $n$  disks—the running time of our algorithm. Our vacuous base case implies that  $T(0) = 0$ , and the more general recursive algorithm implies that  $T(n) = 2T(n - 1) + 1$  for any  $n \geq 1$ . The annihilator method (or guessing and checking by induction) quickly gives us the closed form solution  $T(n) = 2^n - 1$ . In particular, moving a tower of 64 disks requires  $2^{64} - 1 = 18,446,744,073,709,551,615$  individual moves. Thus, even at the impressive rate of one move per second, the monks at Benares will be at work for approximately 585 billion years before tower, temple, and Brahmins alike will crumble into dust, and with a thunderclap the world will vanish.

### 1.3 MergeSort

Mergesort is one of the earliest algorithms proposed for sorting. According to Donald Knuth, it was suggested by John von Neumann as early as 1945.

1. Divide the array  $A[1..n]$  into two subarrays  $A[1..m]$  and  $A[m + 1..n]$ , where  $m = \lfloor n/2 \rfloor$ .
2. Recursively mergesort the subarrays  $A[1..m]$  and  $A[m + 1..n]$ .
3. Merge the newly-sorted subarrays  $A[1..m]$  and  $A[m + 1..n]$  into a single sorted list.

Input:	S	O	R	T	I	N	G	E	X	A	M	P	L
Divide:	S	O	R	T	I	N	G	E	X	A	M	P	L
Recurse:	I	N	O	S	R	T	A	E	G	L	M	P	X
Merge:	A	E	G	I	L	M	N	O	P	S	R	T	X

A Mergesort example.

The first step is completely trivial—we only need to compute the median index  $m$ —and we can delegate the second step to the Recursion Fairy. All the real work is done in the final step; the two sorted subarrays  $A[1..m]$  and  $A[m+1..n]$  can be merged using a simple linear-time algorithm. Here's a complete specification of the Mergesort algorithm; for simplicity, we separate out the merge step as a subroutine.

```
MERGESORT( $A[1..n]$ ):
  if ( $n > 1$ )
     $m \leftarrow \lfloor n/2 \rfloor$ 
    MERGESORT( $A[1..m]$ )
    MERGESORT( $A[m+1..n]$ )
    MERGE( $A[1..n], m$ )
```

```
MERGE( $A[1..n], m$ ):
   $i \leftarrow 1; j \leftarrow m + 1$ 
  for  $k \leftarrow 1$  to  $n$ 
    if  $j > n$ 
       $B[k] \leftarrow A[i]; i \leftarrow i + 1$ 
    else if  $i > m$ 
       $B[k] \leftarrow A[j]; j \leftarrow j + 1$ 
    else if  $A[i] < A[j]$ 
       $B[k] \leftarrow A[i]; i \leftarrow i + 1$ 
    else
       $B[k] \leftarrow A[j]; j \leftarrow j + 1$ 
  for  $k \leftarrow 1$  to  $n$ 
     $A[k] \leftarrow B[k]$ 
```

To prove that the algorithm is correct, we use our old friend induction. We can prove that `MERGE` is correct using induction on the total size of the two subarrays  $A[i..m]$  and  $A[j..n]$  left to be merged into  $B[k..n]$ . The base case, where at least one subarray is empty, is straightforward; the algorithm just copies it into  $B$ . Otherwise, the smallest remaining element is either  $A[i]$  or  $A[j]$ , since both subarrays are sorted, so  $B[k]$  is assigned correctly. The remaining subarrays—either  $A[i+1..m]$  and  $A[j..n]$ , or  $A[i..m]$  and  $A[j+1..n]$ —are merged correctly into  $B[k+1..n]$  by the inductive hypothesis.<sup>4</sup> This completes the proof.

Now we can prove `MERGESORT` correct by another round of straightforward induction. The base cases  $n \leq 1$  are trivial. Otherwise, by the inductive hypothesis, the two smaller subarrays  $A[1..m]$  and  $A[m+1..n]$  are sorted correctly, and by our earlier argument, merged into the correct sorted output.

What's the running time? Since we have a recursive algorithm, we're going to get a recurrence of some sort. `MERGE` clearly takes linear time, since it's a simple for-loop with constant work per iteration. We get the following recurrence for `MERGESORT`:

$$T(n) = T(\lceil n/2 \rceil) + T(\lfloor n/2 \rfloor) + O(n).$$

If we strip out the floors and ceilings, we get the simpler recurrence  $T(n) = 2T(n/2) + O(n)$ . The “all levels equal” case of the recursion tree method (or its corollary, the Master Theorem) immediately implies that  $T(n) = O(n \log n)$ .

**Aside: Domain Transformations.** Wait... what? What do you mean ‘immediately’? How can we just ignore the floors and ceilings? We could always check that  $T(n) = O(n \log n)$  using induction, but there is a simple method for solving recurrences like this directly, called *domain transformation*.

First we overestimate the time bound, once by pretending that the two subproblem sizes are equal, and again to eliminate the ceiling:

$$T(n) \leq 2T(\lceil n/2 \rceil) + O(n) \leq 2T(n/2 + 1) + O(n).$$

Now we define a new function  $S(n) = T(n + \alpha)$ , where  $\alpha$  is a constant chosen so that  $S(n)$  satisfies the familiar recurrence  $S(n) \leq 2S(n/2) + O(n)$ . To figure out the appropriate value for  $\alpha$ , we compare

<sup>4</sup>“The inductive hypothesis” is just a technical nickname for our friend the Recursion Fairy.

two versions of the recurrence for  $T(n + \alpha)$ :

$$\begin{aligned} S(n) \leq 2S(n/2) + O(n) &\implies T(n + \alpha) \leq 2T(n/2 + \alpha) + O(n) \\ T(n) \leq 2T(n/2 + 1) + O(n) &\implies T(n + \alpha) \leq 2T((n + \alpha)/2 + 1) + O(n + \alpha) \end{aligned}$$

For these two recurrences to be equal, we need  $n/2 + \alpha = (n + \alpha)/2 + 1$ , which implies that  $\alpha = 2$ . The recursion tree method tells us that  $S(n) = O(n \log n)$ , so

$$T(n) = S(n - 2) = O((n - 2)\log(n - 2)) = O(n \log n).$$

We can use similar domain transformations to remove floors, ceilings, and lower order terms from any recurrence. But now that we realize this, we don't need to bother grinding through the details ever again!

## 1.4 Quicksort

Quicksort was discovered by Tony Hoare in 1962. In this algorithm, the hard work is splitting the array into subsets so that merging the final result is trivial.

1. Choose a *pivot* element from the array.
2. Split the array into three subarrays containing the items less than the pivot, the pivot itself, and the items bigger than the pivot.
3. Recursively quicksort the first and last subarray.

Input:	S	O	R	T	I	N	G	E	X	A	M	P	L
Choose a pivot:	S	O	R	T	I	N	G	E	X	A	M	P	L
Partition:	M	A	E	G	I	L	N	R	X	O	S	P	T
Recurse:	A	E	G	I	L	M	N	O	P	S	R	T	X

A Quicksort example.

Here's a more formal specification of the Quicksort algorithm. The separate PARTITION subroutine takes the original position of the pivot element as input and returns the post-partition pivot position as output.

<b>QUICKSORT(<math>A[1..n]</math>):</b> <pre> if (<math>n &gt; 1</math>)     Choose a pivot element <math>A[p]</math>     <math>k \leftarrow \text{PARTITION}(A, p)</math>     QUICKSORT(<math>A[1..k - 1]</math>)     QUICKSORT(<math>A[k + 1..n]</math>) </pre>	<b>PARTITION(<math>A[1..n], p</math>):</b> <pre> if (<math>p \neq n</math>)     swap <math>A[p] \leftrightarrow A[n]</math> <math>i \leftarrow 0; j \leftarrow n</math> while (<math>i &lt; j</math>)     repeat <math>i \leftarrow i + 1</math> until (<math>i = j</math> or <math>A[i] \geq A[n]</math>)     repeat <math>j \leftarrow j - 1</math> until (<math>i = j</math> or <math>A[j] \leq A[n]</math>)     if (<math>i &lt; j</math>)         swap <math>A[i] \leftrightarrow A[j]</math> if (<math>i \neq n</math>)     swap <math>A[i] \leftrightarrow A[n]</math> return <math>i</math> </pre>
---	--

Just as we did for mergesort, we need two induction proofs to show that QUICKSORT is correct—weak induction to prove that PARTITION correctly partitions the array, and then straightforward strong induction to prove that QUICKSORT correctly sorts assuming PARTITION is correct. I'll leave the gory details as an exercise for the reader.

The analysis is also similar to mergesort. PARTITION runs in  $O(n)$  time:  $j - i = n$  at the beginning,  $j - i = 0$  at the end, and we do a constant amount of work each time we increment  $i$  or decrement  $j$ . For QUICKSORT, we get a recurrence that depends on  $k$ , the rank of the chosen pivot:

$$T(n) = T(k-1) + T(n-k) + O(n)$$

If we could choose the pivot to be the median element of the array  $A$ , we would have  $k = \lceil n/2 \rceil$ , the two subproblems would be as close to the same size as possible, the recurrence would become

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

and we'd have  $T(n) = O(n \log n)$  by the recursion tree method.

In fact, as we will see later, it is possible to locate the median element in an unsorted array in linear time. However, the algorithm is fairly complicated, and the hidden constant in the  $O()$  notation is quite large. So in practice, programmers settle for something simple, like choosing the first or last element of the array. In this case,  $k$  can be anything from 1 to  $n$ , so we have

$$T(n) = \max_{1 \leq k \leq n} (T(k-1) + T(n-k) + O(n))$$

In the worst case, the two subproblems are completely unbalanced—either  $k = 1$  or  $k = n$ —and the recurrence becomes  $T(n) \leq T(n-1) + O(n)$ . The solution is  $T(n) = O(n^2)$ . Another common heuristic is ‘median of three’—choose three elements (usually at the beginning, middle, and end of the array), and take the middle one as the pivot. Although this is better in practice than just choosing one element, we can still have  $k = 2$  or  $k = n-1$  in the worst case. With the median-of-three heuristic, the recurrence becomes  $T(n) \leq T(1) + T(n-2) + O(n)$ , whose solution is still  $T(n) = O(n^2)$ .

Intuitively, the pivot element will ‘usually’ fall somewhere in the middle of the array, say between  $n/10$  and  $9n/10$ . This suggests that the *average-case* running time is  $O(n \log n)$ . Although this intuition is correct, we are still far from a *proof* that quicksort is usually efficient. We will formalize this intuition about average-case behavior in a later lecture.

## 1.5 The Pattern

Both mergesort and quicksort follow the same general three-step pattern of all divide and conquer algorithms:

1. **Divide** the problem into several *smaller independent* subproblems.
2. **Delegate** each subproblem to the Recursion Fairy to get a sub-solution.
3. **Combine** the sub-solutions together into the final solution.

If the size of any subproblem falls below some constant threshold, the recursion bottoms out. Hopefully, at that point, the problem is trivial, but if not, we switch to a different algorithm instead.

Proving a divide-and-conquer algorithm correct usually involves strong induction. Analyzing the running time requires setting up and solving a recurrence, which often (but unfortunately not always!) can be solved using recursion trees (or, if you insist, the Master Theorem), perhaps after a simple domain transformation.

## 1.6 Median Selection

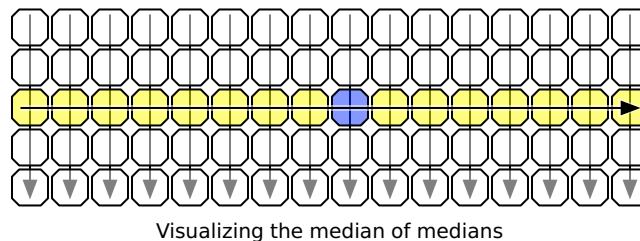
So how do we find the median element of an array in linear time? The following algorithm was discovered by Manuel Blum, Bob Floyd, Vaughan Pratt, Ron Rivest, and Bob Tarjan in the early 1970s. Their algorithm actually solves the more general problem of selecting the  $k$ th largest element in an array, using the following recursive divide-and-conquer strategy. The subroutine PARTITION is the same as the one used in QUICKSORT.

```
SELECT( $A[1..n], k$ ):
  if  $n \leq 25$ 
    use brute force
  else
     $m \leftarrow \lceil n/5 \rceil$ 
    for  $i \leftarrow 1$  to  $m$ 
       $B[i] \leftarrow \text{SELECT}(A[5i - 4..5i], 3)$     «Brute force!»
     $mom \leftarrow \text{SELECT}(B[1..m], \lfloor m/2 \rfloor)$     «Recursion!»
     $r \leftarrow \text{PARTITION}(A[1..n], mom)$ 
    if  $k < r$ 
      return  $\text{SELECT}(A[1..r - 1], k)$     «Recursion!»
    else if  $k > r$ 
      return  $\text{SELECT}(A[r + 1..n], k - r)$     «Recursion!»
    else
      return  $mom$ 
```

If the input array is too large to handle by brute force, we divide it into  $\lceil n/5 \rceil$  blocks, each containing exactly 5 elements, except possibly the last. (If the last block isn't full, just throw in a few  $\infty$ s.) We find the median of each block by brute force and collect those medians into a new array. Then we recursively compute the median of the new array (the median of medians — hence ‘mom’) and use it to partition the input array. Finally, either we get lucky and the median-of-medians is the  $k$ th largest element of  $A$ , or we recursively search one of the two subarrays.

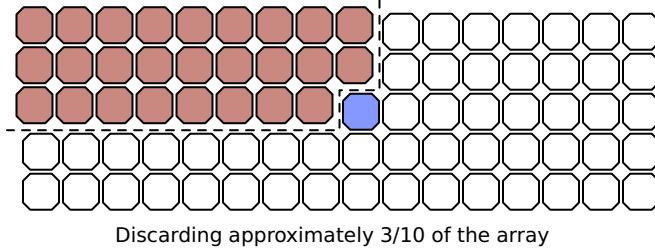
The key insight is that these two subarrays cannot be too large or too small. The median-of-medians is larger than  $\lceil \lceil n/5 \rceil / 2 \rceil - 1 \approx n/10$  medians, and each of those medians is larger than two other elements in its block. In other words, the median-of-medians is larger than at least  $3n/10$  elements in the input array. Symmetrically,  $mom$  is smaller than at least  $3n/10$  input elements. Thus, in the worst case, the final recursive call searches an array of size  $7n/10$ .

We can visualize the algorithm's behavior by drawing the input array as a  $5 \times \lceil n/5 \rceil$  grid, which each column represents five consecutive elements. For purposes of illustration, imagine that we sort every column from top down, and then we sort the columns by their middle element. (Let me emphasize that the *algorithm* doesn't actually do this!) In this arrangement, the median-of-medians is the element closest to the center of the grid.



The left half of the first three rows of the grid contains  $3n/10$  elements, each of which is smaller than the median-of-medians. If the element we're looking for is larger than the median-of-medians,

our algorithm will throw away *everything* smaller than the median-of-median, including those  $3n/10$  elements, before recursing. A symmetric argument applies when our target element is smaller than the median-of-medians.



We conclude that the worst-case running time of the algorithm obeys the following recurrence:

$$T(n) \leq O(n) + T(n/5) + T(7n/10).$$

The recursion tree method implies the solution  $T(n) = O(n)$ .

Finer analysis reveals that the hidden constants are quite large, even if we count only comparisons; this is not a practical algorithm for small inputs. (In particular, mergesort uses fewer comparisons in the worst case when  $n < 4,000,000$ .) Selecting the median of 5 elements requires **at most 6 comparisons**, so we need at most  $6n/5$  comparisons to set up the recursive subproblem. We need another  $n - 1$  comparisons to partition the array after the recursive call returns. So the actual recurrence is

$$T(n) \leq 11n/5 + T(n/5) + T(7n/10).$$

The recursion tree method implies the upper bound

$$T(n) \leq \frac{11n}{5} \sum_{i \geq 0} \left(\frac{9}{10}\right)^i = \frac{11n}{5} \cdot 10 = 22n.$$

## 1.7 Multiplication

Adding two  $n$ -digit numbers takes  $O(n)$  time by the standard iterative ‘ripple-carry’ algorithm, using a lookup table for each one-digit addition. Similarly, multiplying an  $n$ -digit number by a one-digit number takes  $O(n)$  time, using essentially the same algorithm.

What about multiplying two  $n$ -digit numbers? At least in the United States, every grade school student (supposedly) learns to multiply by breaking the problem into  $n$  one-digit multiplications and  $n$  additions:

$$\begin{array}{r}
 31415962 \\
 \times 27182818 \\
 \hline
 251327696 \\
 31415962 \\
 251327696 \\
 62831924 \\
 251327696 \\
 31415962 \\
 219911734 \\
 62831924 \\
 \hline
 853974377340916
 \end{array}$$

We could easily formalize this algorithm as a pair of nested for-loops. The algorithm runs in  $O(n^2)$  time—altogether, there are  $O(n^2)$  digits in the partial products, and for each digit, we spend constant time. The Egyptian/Russian peasant multiplication algorithm described in the first lecture also runs in  $O(n^2)$  time.

We can get a more efficient algorithm by exploiting the following identity:

$$(10^m a + b)(10^m c + d) = 10^{2m} ac + 10^m(bc + ad) + bd$$

Here is a divide-and-conquer algorithm that computes the product of two  $n$ -digit numbers  $x$  and  $y$ , based on this formula. Each of the four sub-products  $e, f, g, h$  is computed recursively. The last line does not involve any multiplications, however; to multiply by a power of ten, we just shift the digits and fill in the right number of zeros.

```
MULTIPLY( $x, y, n$ ):
  if  $n = 1$ 
    return  $x \cdot y$ 
  else
     $m \leftarrow \lceil n/2 \rceil$ 
     $a \leftarrow \lfloor x/10^m \rfloor$ ;  $b \leftarrow x \bmod 10^m$ 
     $d \leftarrow \lfloor y/10^m \rfloor$ ;  $c \leftarrow y \bmod 10^m$ 
     $e \leftarrow \text{MULTIPLY}(a, c, m)$ 
     $f \leftarrow \text{MULTIPLY}(b, d, m)$ 
     $g \leftarrow \text{MULTIPLY}(b, c, m)$ 
     $h \leftarrow \text{MULTIPLY}(a, d, m)$ 
    return  $10^{2m}e + 10^m(g + h) + f$ 
```

You can easily prove by induction that this algorithm is correct. The running time for this algorithm is given by the recurrence

$$T(n) = 4T(\lceil n/2 \rceil) + O(n), \quad T(1) = 1,$$

which solves to  $T(n) = O(n^2)$  by the recursion tree method (after a simple domain transformation). Hmm... I guess this didn't help after all.

But there's a trick, first published by Anatoliï Karatsuba in 1962.<sup>5</sup> We can compute the middle coefficient  $bc + ad$  using only *one* recursive multiplication, by exploiting yet another bit of algebra:

$$ac + bd - (a - b)(c - d) = bc + ad$$

This trick lets us replace the last three lines in the previous algorithm as follows:

```
FASTMULTIPLY( $x, y, n$ ):
  if  $n = 1$ 
    return  $x \cdot y$ 
  else
     $m \leftarrow \lceil n/2 \rceil$ 
     $a \leftarrow \lfloor x/10^m \rfloor$ ;  $b \leftarrow x \bmod 10^m$ 
     $d \leftarrow \lfloor y/10^m \rfloor$ ;  $c \leftarrow y \bmod 10^m$ 
     $e \leftarrow \text{FASTMULTIPLY}(a, c, m)$ 
     $f \leftarrow \text{FASTMULTIPLY}(b, d, m)$ 
     $g \leftarrow \text{FASTMULTIPLY}(a - b, c - d, m)$ 
    return  $10^{2m}e + 10^m(e + f - g) + f$ 
```

<sup>5</sup>However, the same basic trick was used non-recursively by Gauss in the 1800s to multiply complex numbers using only three real multiplications.

The running time of Karatsuba's FASTMULTIPLY algorithm is given by the recurrence

$$T(n) \leq 3T(\lceil n/2 \rceil) + O(n), \quad T(1) = 1.$$

After a domain transformation, we can plug this into a recursion tree to get the solution  $T(n) = O(n^{\lg 3}) = O(n^{1.585})$ , a significant improvement over our earlier quadratic-time algorithm.<sup>6</sup>

Of course, in practice, all this is done in binary instead of decimal.

We can take this idea even further, splitting the numbers into more pieces and combining them in more complicated ways, to get even faster multiplication algorithms. Ultimately, this idea leads to the development of the *Fast Fourier transform*, a more complicated divide-and-conquer algorithm that can be used to multiply two  $n$ -digit numbers in  $O(n \log n)$  time.<sup>7</sup> We'll talk about Fast Fourier transforms in detail in the next lecture.

## 1.8 Exponentiation

Given a number  $a$  and a positive integer  $n$ , suppose we want to compute  $a^n$ . The standard naïve method is a simple for-loop that does  $n - 1$  multiplications by  $a$ :

```
SLOWPOWER( $a, n$ ):
   $x \leftarrow a$ 
  for  $i \leftarrow 2$  to  $n$ 
     $x \leftarrow x \cdot a$ 
  return  $x$ 
```

This iterative algorithm requires  $n$  multiplications.

Notice that the input  $a$  could be an integer, or a rational, or a floating point number. In fact, it doesn't need to be a number at all, as long as it's something that we know how to multiply. For example, the same algorithm can be used to compute powers modulo some finite number (an operation commonly used in cryptography algorithms) or to compute powers of matrices (an operation used to evaluate recurrences and to compute shortest paths in graphs). All we really require is that  $a$  belong to a multiplicative group.<sup>8</sup> Since we don't know what kind of things we're multiplying, we can't know how long a multiplication takes, so we're forced analyze the running time in terms of the number of multiplications.

There is a much faster divide-and-conquer method, using the following simple recursive formula:

$$a^n = a^{\lfloor n/2 \rfloor} \cdot a^{\lceil n/2 \rceil}.$$

What makes this approach more efficient is that once we compute the first factor  $a^{\lfloor n/2 \rfloor}$ , we can compute the second factor  $a^{\lceil n/2 \rceil}$  using at most one more multiplication.

<sup>6</sup>Karatsuba actually proposed an algorithm based on the formula  $(a+c)(b+d) - ac - bd = bc + ad$ . This algorithm also runs in  $O(n^{\lg 3})$  time, but the actual recurrence is a bit messier:  $a - b$  and  $c - d$  are still  $m$ -digit numbers, but  $a + b$  and  $c + d$  might have  $m + 1$  digits. The simplification presented here is due to Donald Knuth.

<sup>7</sup>This fast algorithm for multiplying integers using FFTs was discovered by Arnold Schönhage and Volker Strassen in 1971. The  $O(n \log n)$  running time requires the standard assumption that  $O(\log n)$ -bit integer arithmetic can be performed in constant time; the number of bit operations is  $O(n \log n \log \log n)$ .

<sup>8</sup>A *multiplicative group*  $(G, \otimes)$  is a set  $G$  and a function  $\otimes : G \times G \rightarrow G$ , satisfying three axioms:

1. There is a *unit* element  $1 \in G$  such that  $1 \otimes g = g \otimes 1$  for any element  $g \in G$ .
2. Any element  $g \in G$  has an *inverse* element  $g^{-1} \in G$  such that  $g \otimes g^{-1} = g^{-1} \otimes g = 1$ .
3. The function is *associative*: for any elements  $f, g, h \in G$ , we have  $f \otimes (g \otimes h) = (f \otimes g) \otimes h$ .

```

FASTPOWER( $a, n$ ):
  if  $n = 1$ 
    return  $a$ 
  else
     $x \leftarrow \text{FASTPOWER}(a, \lfloor n/2 \rfloor)$ 
    if  $n$  is even
      return  $x \cdot x$ 
    else
      return  $x \cdot x \cdot a$ 

```

The total number of multiplications satisfies the recurrence  $T(n) \leq T(\lfloor n/2 \rfloor) + 2$ , with the base case  $T(1) = 0$ . After a domain transformation, recursion trees give us the solution  $T(n) = O(\log n)$ .

Incidentally, this algorithm is asymptotically optimal—any algorithm for computing  $a^n$  must perform at least  $\Omega(\log n)$  multiplications. In fact, when  $n$  is a power of two, this algorithm is *exactly* optimal. However, there are slightly faster methods for other values of  $n$ . For example, our divide-and-conquer algorithm computes  $a^{15}$  in six multiplications ( $a^{15} = a^7 \cdot a^7 \cdot a$ ;  $a^7 = a^3 \cdot a^3 \cdot a$ ;  $a^3 = a \cdot a \cdot a$ ), but only five multiplications are necessary ( $a \rightarrow a^2 \rightarrow a^3 \rightarrow a^5 \rightarrow a^{10} \rightarrow a^{15}$ ). Nobody knows of an efficient algorithm that always uses the minimum possible number of multiplications.<sup>9</sup>

## Exercises

1. (a) Professor George O'Jungle has a 27-node binary tree, in which every node is labeled with a unique letter of the Roman alphabet or the character &. Preorder and postorder traversals of the tree visit the nodes in the following order:
    - Preorder: **I Q J H L E M V O T S B R G Y Z K C A & F P N U D W X**
    - Postorder: **H E M L J V Q S G Y R Z B T C P U D N F W & X A K O I**
 Draw George's binary tree.
  - (b) Describe and analyze a recursive algorithm for reconstructing a binary tree, given its preorder and postorder node sequences.
  - (c) Describe and analyze a recursive algorithm for reconstructing a binary tree, given its preorder and *inorder* node sequences.
- 
2. Consider a  $2^n \times 2^n$  chessboard with one (arbitrarily chosen) square removed.
    - (a) Prove that any such chessboard can be tiled without gaps or overlaps by L-shaped pieces, each composed of 3 squares.
    - (b) Describe and analyze an algorithm to compute such a tiling, given the integer  $n$  and two  $n$ -bit integers representing the row and column of the missing square. The output is a list of the positions and orientations of  $(4^n - 1)/3$  tiles. Your algorithm should run in  $O(4^n)$  time.
  3. Prove that the original recursive Tower of Hanoi algorithm is *exactly equivalent* to each of the following non-recursive algorithms; in other words, prove that all three algorithms move the same disks, to and from the same needles, in the same order. The needles are labeled 0, 1, and 2, and our problem is to move a stack of  $n$  disks from needle 0 to needle 2 (as shown on page ??).

---

<sup>9</sup>On the other hand, nobody knows whether finding the minimum number of multiplications is NP-hard, either.

(a) Follow these four rules:

- Never move the same disk twice in a row.
- If  $n$  is even, always move the smallest disk forward ( $\dots \rightarrow 0 \rightarrow 1 \rightarrow 2 \rightarrow 0 \rightarrow \dots$ ).
- If  $n$  is odd, always move the smallest disk backward ( $\dots \rightarrow 0 \rightarrow 2 \rightarrow 1 \rightarrow 0 \rightarrow \dots$ ).
- When there is no move that satisfies the other rules, the puzzle is solved.

(b) Let  $\rho(n)$  denote the smallest integer  $k$  such that  $n/2^k$  is not an integer. For example,  $\rho(42) = 2$ , because  $42/2^1$  is an integer but  $42/2^2$  is not. (Equivalently,  $\rho(n)$  is one more than the position of the least significant 1 in the binary representation of  $n$ .) The function  $\rho(n)$  is sometimes called the ‘ruler’ function, because its behavior resembles the marks on a ruler:

1, 2, 1, 3, 1, 2, 1, 4, 1, 2, 1, 3, 1, 2, 1, 5, 1, 2, 1, 3, 1, 2, 1, 4, 1, 2, 1, 3, 1, 2, 1, 6, 1, 2, 1, 3, 1, ...

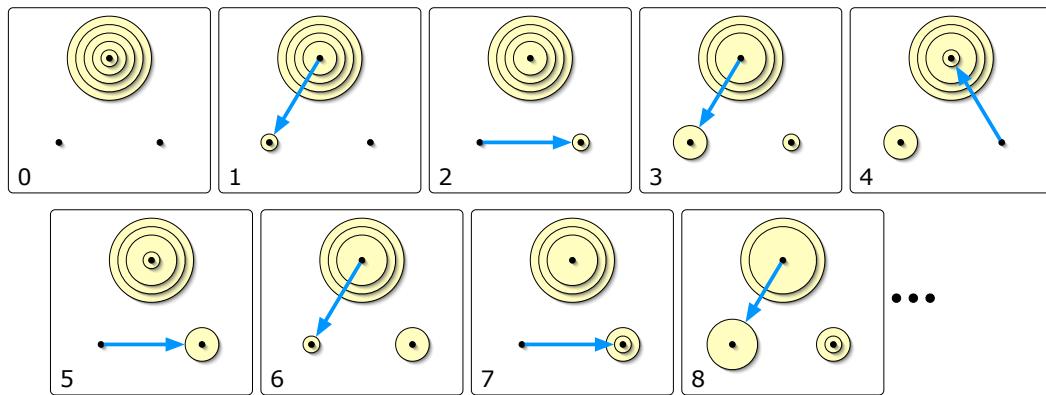
Here’s the non-recursive algorithm in one line:

**In step  $i$ , move disk  $\rho(i)$  forward if  $n - i$  is even, backward if  $n - i$  is odd.**

When this rule requires us to move disk  $n + 1$ , the algorithm ends.

4. Consider the following restricted variants of the Tower of Hanoi puzzle. In each problem, the needles are numbered 0, 1, and 2, as in problem ??, and your task is to move a stack of  $n$  disks from needle 1 to needle 2.

- Suppose you are forbidden to move any disk directly between needle 1 and needle 2; every move must involve needle 0. Describe an algorithm to solve this version of the puzzle in as few moves as possible. *Exactly* how many moves does your algorithm make?
- Suppose you are only allowed to move disks from needle 0 to needle 2, from needle 2 to needle 1, or from needle 1 to needle 0. Equivalently, Suppose the needles are arranged in a circle and numbered in clockwise order, and you are only allowed to move disks counterclockwise. Describe an algorithm to solve this version of the puzzle in as few moves as possible. *Exactly* how many moves does your algorithm make?



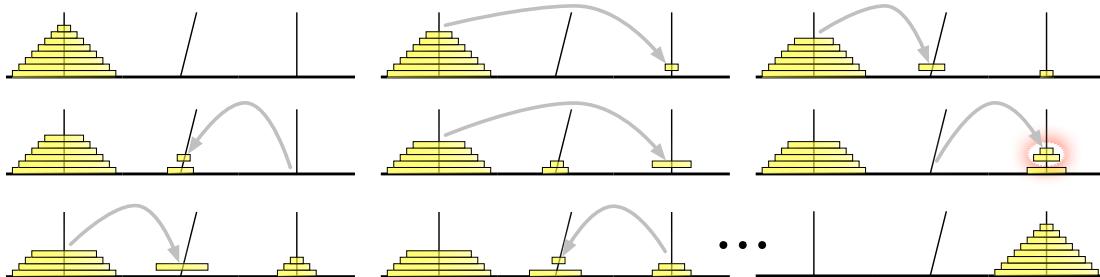
A top view of the first eight moves in a clockwise Towers of Hanoi solution

- ★(c) Finally, suppose your only restriction is that you may never move a disk directly from needle 1 to needle 2. Describe an algorithm to solve this version of the puzzle in as few moves as possible. How many moves does your algorithm make? [Hint: This is considerably harder to analyze than the other two variants.]

5. A less familiar chapter in the Tower of Hanoi's history is its brief relocation of the temple from Benares to Pisa in the early 13th century. The relocation was organized by the wealthy merchant-mathematician Leonardo Fibonacci, at the request of the Holy Roman Emperor Frederick II, who had heard reports of the temple from soldiers returning from the Crusades. The Towers of Pisa and their attendant monks became famous, helping to establish Pisa as a dominant trading center on the Italian peninsula.

Unfortunately, almost as soon as the temple was moved, one of the diamond needles began to lean to one side. To avoid the possibility of the leaning tower falling over from too much use, Fibonacci convinced the priests to adopt a more relaxed rule: *Any number of disks on the leaning needle can be moved together to another needle in a single move*. It was still forbidden to place a larger disk on top of a smaller disk, and disks had to be moved one at a time *onto* the leaning needle or between the two vertical needles.

Thanks to Fibonacci's new rule, the priests could bring about the end of the universe somewhat faster from Pisa than they could than could from Benares. Fortunately, the temple was moved from Pisa back to Benares after the newly crowned Pope Gregory IX excommunicated Frederick II, making the local priests less sympathetic to hosting foreign heretics with strange mathematical habits. Soon afterward, a bell tower was erected on the spot where the temple once stood; it too began to lean almost immediately.



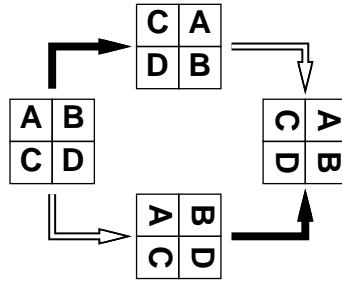
The Towers of Pisa. In the fifth move, two disks are taken off the leaning needle.

Describe an algorithm to transfer a stack of  $n$  disks from one *vertical* needle to the other *vertical* needle, using the smallest possible number of moves. *Exactly* how many moves does your algorithm perform?

6. Most graphics hardware includes support for a low-level operation called *blit*, or **block transfer**, which quickly copies a rectangular chunk of a pixel map (a two-dimensional array of pixel values) from one location to another. This is a two-dimensional version of the standard C library function `memcpy()`.

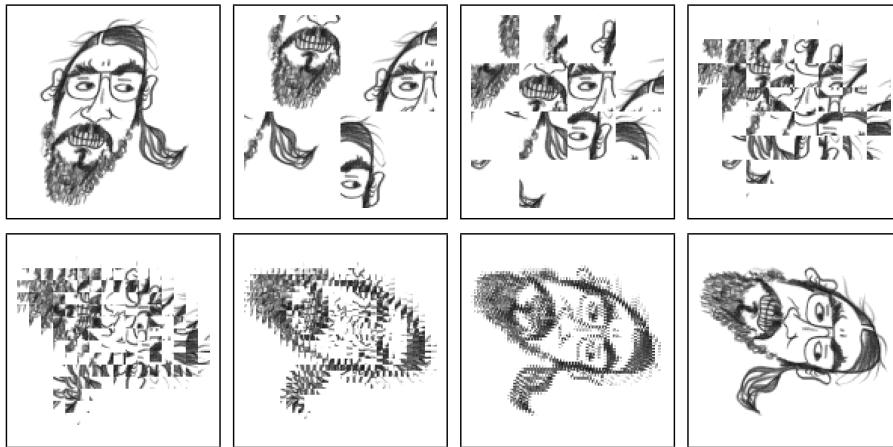
Suppose we want to rotate an  $n \times n$  pixel map  $90^\circ$  clockwise. One way to do this, at least when  $n$  is a power of two, is to split the pixel map into four  $n/2 \times n/2$  blocks, move each block to its proper position using a sequence of five blits, and then recursively rotate each block. Alternately, we could *first* recursively rotate the blocks and *then* blit them into place.

- (a) Prove that both versions of the algorithm are correct when  $n$  is a power of two.
- (b) *Exactly* how many blits does the algorithm perform when  $n$  is a power of two?
- (c) Describe how to modify the algorithm so that it works for arbitrary  $n$ , not just powers of two. How many blits does your modified algorithm perform?



Two algorithms for rotating a pixel map.

Black arrows indicate blitting the blocks into place; white arrows indicate recursively rotating the blocks.



The first rotation algorithm (blit then recurse) in action.

- (d) What is your algorithm's running time if a  $k \times k$  blit takes  $O(k^2)$  time?
- (e) What if a  $k \times k$  blit takes only  $O(k)$  time?
7. An *inversion* in an array  $A[1..n]$  is a pair of indices  $(i, j)$  such that  $i < j$  and  $A[i] > A[j]$ . The number of inversions in an  $n$ -element array is between 0 (if the array is sorted) and  $\binom{n}{2}$  (if the array is sorted backward).

Describe and analyze an algorithm to count the number of inversions in an  $n$ -element array in  $O(n \log n)$  time. [Hint: Modify mergesort.]

8. (a) Prove that the following algorithm actually sorts its input.

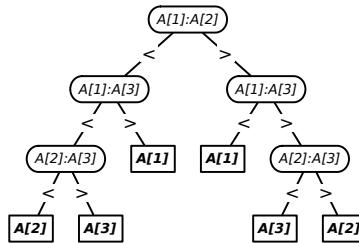
```

STOOGESORT(A[0..n - 1]) :
  if n = 2 and A[0] > A[1]
    swap A[0] ↔ A[1]
  else if n > 2
    m = ⌈2n/3⌉
    STOOGESORT(A[0..m - 1])
    STOOGESORT(A[n - m..n - 1])
    STOOGESORT(A[0..m - 1])
  
```

- (b) Would STOOGESORT still sort correctly if we replaced  $m = \lceil 2n/3 \rceil$  with  $m = \lfloor 2n/3 \rfloor$ ? Justify your answer.

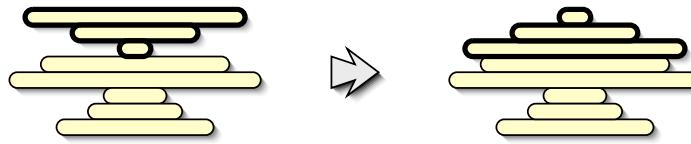
- (c) State a recurrence (including the base case(s)) for the number of comparisons executed by STOOGESORT.
- (d) Solve the recurrence, and prove that your solution is correct. [*Hint: Ignore the ceiling.*]
- (e) Prove that the number of swaps executed by STOOGESORT is at most  $\binom{n}{2}$ .
9. Describe an algorithm to compute the median of an array  $A[1..5]$  of distinct numbers using at most 6 comparisons.

Instead of writing pseudocode, describe your algorithm using a decision tree: A binary tree where each internal node is a comparison of the form “ $A[i] \geq A[j]?$ ” and each leaf is an index into the array.



Finding the median of a 3-element array using at most 3 comparisons

10. Suppose you are given a stack of  $n$  pancakes of different sizes. You want to sort the pancakes so that smaller pancakes are on top of larger pancakes. The only operation you can perform is a *flip*—insert a spatula under the top  $k$  pancakes, for some integer  $k$  between 1 and  $n$ , and flip them all over.



Flipping the top three pancakes.

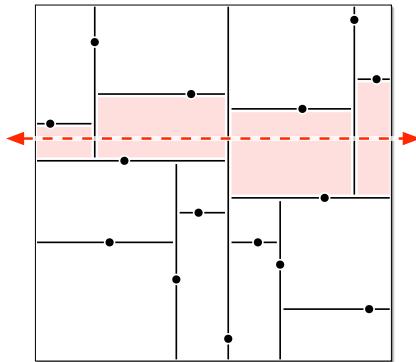
- (a) Describe an algorithm to sort an arbitrary stack of  $n$  pancakes using as few flips as possible. *Exactly* how many flips does your algorithm perform in the worst case?
- (b) Now suppose one side of each pancake is burned. Describe an algorithm to sort an arbitrary stack of  $n$  pancakes, so that the burned side of every pancake is facing down, using as few flips as possible. *Exactly* how many flips does your algorithm perform in the worst case?
11. Suppose we are given an array  $A[1..n]$  with the special property that  $A[1] \geq A[2]$  and  $A[n-1] \leq A[n]$ . We say that an element  $A[x]$  is a *local minimum* if it is less than or equal to both its neighbors, or more formally, if  $A[x-1] \geq A[x]$  and  $A[x] \leq A[x+1]$ . For example, there are six local minima in the following array:

9	7	7	2	1	3	7	5	4	7	3	3	4	8	6	9
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

We can obviously find a local minimum in  $O(n)$  time by scanning through the array. Describe and analyze an algorithm that finds a local minimum in  $O(\log n)$  time. [*Hint: With the given boundary conditions, the array **must** have at least one local minimum. Why?*]

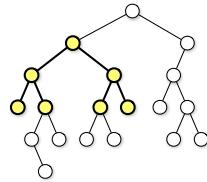
12. You are a contestant on the hit game show “Beat Your Neighbors!” You are presented with an  $m \times n$  grid of boxes, each containing a unique number. It costs \$100 to open a box. Your goal is to find a box whose number is larger than its neighbors in the grid (above, below, left, and right). If you spend less money than any of your opponents, you win a week-long trip for two to Las Vegas and a year’s supply of Rice-A-Roni™, to which you are hopelessly addicted.
- (a) Suppose  $m = 1$ . Describe an algorithm that finds a number that is bigger than any of its neighbors. How many boxes does your algorithm open in the worst case?
  - \*(b) Suppose  $m = n$ . Describe an algorithm that finds a number that is bigger than any of its neighbors. How many boxes does your algorithm open in the worst case?
  - \*(c) Prove that your solution to part (b) is optimal up to a constant factor.
13. (a) Suppose we are given two sorted arrays  $A[1..n]$  and  $B[1..n]$  and an integer  $k$ . Describe an algorithm to find the  $k$ th smallest element in the union of  $A$  and  $B$  in  $\Theta(\log n)$  time. For example, if  $k = 1$ , your algorithm should return the smallest element of  $A \cup B$ ; if  $k = n$ , your algorithm should return the median of  $A \cup B$ .) You can assume that the arrays contain no duplicate elements. [*Hint: First solve the special case  $k = n$ .*]
- (b) Now suppose we are given *three* sorted arrays  $A[1..n]$ ,  $B[1..n]$ , and  $C[1..n]$ , and an integer  $k$ . Describe an algorithm to find the  $k$ th smallest element in  $A \cup B \cup C$  in  $O(\log n)$  time.
  - (c) Finally, suppose we are given a two dimensional array  $A[1..m][1..n]$  in which every row  $A[i][\cdot]$  is sorted, and an integer  $k$ . Describe an algorithm to find the  $k$ th smallest element in  $A$  as quickly as possible. How does the running time of your algorithm depend on  $m$ ? [*Hint: Use the linear-time SELECT algorithm as a subroutine.*]
14. (a) Describe an algorithm that sorts an input array  $A[1..n]$  by calling a subroutine  $SQRT SORT(k)$ , which sorts the subarray  $A[k + 1..k + \sqrt{n}]$  in place, given an arbitrary integer  $k$  between 0 and  $n - \sqrt{n}$  as input. (To simplify the problem, assume that  $\sqrt{n}$  is an integer.) Your algorithm is **only** allowed to inspect or modify the input array by calling  $SQRT SORT$ ; in particular, your algorithm must not directly compare, move, or copy array elements. How many times does your algorithm call  $SQRT SORT$  in the worst case?
- (b) Prove that your algorithm from part (a) is optimal up to constant factors. In other words, if  $f(n)$  is the number of times your algorithm calls  $SQRT SORT$ , prove that no algorithm can sort using  $o(f(n))$  calls to  $SQRT SORT$ . [*Hint: See Lecture 19.*]
  - (c) Now suppose  $SQRT SORT$  is implemented recursively, by calling your sorting algorithm from part (a). For example, at the second level of recursion, the algorithm is sorting arrays roughly of size  $n^{1/4}$ . What is the worst-case running time of the resulting sorting algorithm? (To simplify the analysis, assume that the array size  $n$  has the form  $2^{2^k}$ , so that repeated square roots are always integers.)
15. Suppose we have  $n$  points scattered inside a two-dimensional box. A *kd-tree* recursively subdivides the points as follows. First we split the box into two smaller boxes with a *vertical* line, then we split each of those boxes with *horizontal* lines, and so on, always alternating between horizontal and vertical splits. Each time we split a box, the splitting line partitions the rest of the interior points as *evenly as possible* by passing through a median point inside the box (*not* on its boundary).

If a box doesn't contain any points, we don't split it any more; these final empty boxes are called *cells*.



A kd-tree for 15 points. The dashed line crosses the four shaded cells.

- (a) How many cells are there, as a function of  $n$ ? Prove your answer is correct.
  - (b) In the worst case, *exactly* how many cells can a horizontal line cross, as a function of  $n$ ? Prove your answer is correct. Assume that  $n = 2^k - 1$  for some integer  $k$ . [Hint: There is more than one function  $f$  such that  $f(16) = 4$ .]
  - (c) Suppose we have  $n$  points stored in a kd-tree. Describe and analyze an algorithm that counts the number of points above a horizontal line (such as the dashed line in the figure) as quickly as possible. [Hint: Use part (b).]
  - (d) Describe an analyze an efficient algorithm that counts, given a kd-tree storing  $n$  points, the number of points that lie inside a rectangle  $R$  with horizontal and vertical sides. [Hint: Use part (c).]
16. You are at a political convention with  $n$  delegates, each one a member of exactly one political party. It is impossible to tell which political party any delegate belongs to; in particular, you will be summarily ejected from the convention if you ask. However, you can determine whether any two delegates belong to the *same* party or not by introducing them to each other—members of the same party always greet each other with smiles and friendly handshakes; members of different parties always greet each other with angry stares and insults.
- (a) Suppose a majority (more than half) of the delegates are from the same political party. Describe an efficient algorithm that identifies a member (*any* member) of the majority party.
  - (b) Now suppose exactly  $k$  political parties are represented at the convention and one party has a *plurality*: more delegates belong to that party than to any other. Present a practical procedure to pick a person from the plurality political party as parsimoniously as possible. (Please.)
17. For this problem, a *subtree* of a binary tree means any connected subgraph. A binary tree is *complete* if every internal node has two children, and every leaf has exactly the same depth. Describe and analyze a recursive algorithm to compute the *largest complete subtree* of a given binary tree. Your algorithm should return the root and the depth of this subtree.



The largest complete subtree of this binary tree has depth 2.

18. Consider the following classical recursive algorithm for computing the factorial  $n!$  of a non-negative integer  $n$ :

```
FACTORIAL( $n$ ):
  if  $n = 0$ 
    return 0
  else
    return  $n \cdot \text{FACTORIZATION}(n - 1)$ 
```

- (a) How many multiplications does this algorithm perform?
- (b) How many bits are required to write  $n!$  in binary? Express your answer in the form  $\Theta(f(n))$ , for some familiar function  $f(n)$ . [Hint: Use Stirling's approximation:  $n! \approx \sqrt{2\pi n} \cdot (n/e)^n$ .]
- (c) Your answer to (b) should convince you that the number of multiplications is *not* a good estimate of the actual running time of FACTORIAL. The grade-school multiplication algorithm takes  $O(k \cdot l)$  time to multiply a  $k$ -digit number and an  $l$ -digit number. What is the running time of FACTORIAL if we use this multiplication algorithm as a subroutine?
- (d) The following algorithm also computes  $n!$ , but groups the multiplication differently:

```
FACTORIZATION2( $n, m$ ):      «Compute  $n!/(n - m)!$ »
  if  $m = 0$ 
    return 1
  else if  $m = 1$ 
    return  $n$ 
  else
    return FACTORIZATION2( $n, \lfloor m/2 \rfloor$ ) · FACTORIZATION2( $n - \lfloor m/2 \rfloor, \lceil m/2 \rceil$ )
```

What is the running time of FACTORIZATION2( $n, n$ ) if we use grade-school multiplication? [Hint: Ignore the floors and ceilings.]

- (e) Describe and analyze a variant of Karatsuba's algorithm that can multiply any  $k$ -digit number and any  $l$ -digit number, where  $k \geq l$ , in  $O(k \cdot l^{\lg 3 - 1}) = O(k \cdot l^{0.585})$  time.
- (f) What are the running times of FACTORIAL( $n$ ) and FACTORIZATION2( $n, n$ ) if we use the modified Karatsuba multiplication from part (e)?

**Calvin:** Here's another math problem I can't figure out. What's  $9+4$ ?

**Hobbes:** Ooh, that's a tricky one. You have to use calculus and imaginary numbers for this.

**Calvin:** IMAGINARY NUMBERS?!

**Hobbes:** You know, eleventeen, thirty-twelve, and all those. It's a little confusing at first.

**Calvin:** How did YOU learn all this? You've never even gone to school!

**Hobbes:** Instinct. Tigers are born with it.

— “Calvin and Hobbes” (January 6, 1998)

*It needs evaluation*

*So let the games begin*

*A heinous crime, a show of force*

*A murder would be nice, of course*

— “Bad Horse’s Letter”, Dr. Horrible’s Sing-Along Blog (2008)

## \*2 Fast Fourier Transforms

### 2.1 Polynomials

In this lecture we'll talk about algorithms for manipulating *polynomials*: functions of one variable built from additions subtractions, and multiplications (but no divisions). The most common representation for a polynomial  $p(x)$  is as a sum of weighted powers of a variable  $x$ :

$$p(x) = \sum_{j=0}^n a_j x^j.$$

The numbers  $a_j$  are called *coefficients*. The *degree* of the polynomial is the largest power of  $x$ ; in the example above, the degree is  $n$ . Any polynomial of degree  $n$  can be specified by a sequence of  $n + 1$  coefficients. Some of these coefficients may be zero, but not the  $n$ th coefficient, because otherwise the degree would be less than  $n$ .

Here are three of the most common operations that are performed with polynomials:

- **Evaluate:** Give a polynomial  $p$  and a number  $x$ , compute the number  $p(x)$ .
- **Add:** Give two polynomials  $p$  and  $q$ , compute a polynomial  $r = p + q$ , so that  $r(x) = p(x) + q(x)$  for all  $x$ . If  $p$  and  $q$  both have degree  $n$ , then their sum  $p + q$  also has degree  $n$ .
- **Multiply:** Give two polynomials  $p$  and  $q$ , compute a polynomial  $r = p \cdot q$ , so that  $r(x) = p(x) \cdot q(x)$  for all  $x$ . If  $p$  and  $q$  both have degree  $n$ , then their product  $p \cdot q$  has degree  $2n$ .

Suppose we represent a polynomial of degree  $n$  as an array of  $n + 1$  coefficients  $P[0..n]$ , where  $P[j]$  is the coefficient of the  $x^j$  term. We learned simple algorithms for all three of these operations in high-school algebra:

```
EVALUATE( $P[0..n], x$ ):
   $X \leftarrow 1$      $\langle X = x^j \rangle$ 
   $y \leftarrow 0$ 
  for  $j \leftarrow 0$  to  $n$ 
     $y \leftarrow y + P[j] \cdot X$ 
     $X \leftarrow X \cdot x$ 
  return  $y$ 
```

```
ADD( $P[0..n], Q[0..n]$ ):
  for  $j \leftarrow 0$  to  $n$ 
     $R[j] \leftarrow P[j] + Q[j]$ 
  return  $R[0..n]$ 
```

```
MULTIPLY( $P[0..n], Q[0..m]$ ):
  for  $j \leftarrow 0$  to  $n + m$ 
     $R[j] \leftarrow 0$ 
  for  $j \leftarrow 0$  to  $n$ 
    for  $k \leftarrow 0$  to  $m$ 
       $R[j+k] \leftarrow R[j+k] + P[j] \cdot Q[k]$ 
  return  $R[0..n+m]$ 
```

EVALUATE uses  $O(n)$  arithmetic operations.<sup>1</sup> This is the best we can hope for, but we can cut the number of multiplications in half using *Horner's rule*:

$$p(x) = a_0 + x(a_1 + x(a_2 + \dots + x a_n)).$$

```

HORNER( $P[0..n]$ ,  $x$ ):
   $y \leftarrow P[n]$ 
  for  $i \leftarrow n - 1$  downto 0
     $y \leftarrow x \cdot y + P[i]$ 
  return  $y$ 
```

The addition algorithm also runs in  $O(n)$  time, and this is clearly the best we can do.

The multiplication algorithm, however, runs in  $O(n^2)$  time. In the previous lecture, we saw a divide and conquer algorithm (due to Karatsuba) for multiplying two  $n$ -bit integers in only  $O(n^{\lg 3})$  steps; precisely the same algorithm can be applied here. Even cleverer divide-and-conquer strategies lead to multiplication algorithms whose running times are arbitrarily close to linear— $O(n^{1+\epsilon})$  for your favorite value  $\epsilon > 0$ —but with great cleverness comes great confusion. These algorithms are difficult to understand, even more difficult to implement correctly, and not worth the trouble in practice thanks to large constant factors.

## 2.2 Alternate Representations

Part of what makes multiplication so much harder than the other two operations is our input representation. Coefficients vectors are the most common representation for polynomials, but there are at least two other useful representations.

### 2.2.1 Roots

The Fundamental Theorem of Algebra states that every polynomial  $p$  of degree  $n$  has exactly  $n$  roots  $r_1, r_2, \dots, r_n$  such that  $p(r_j) = 0$  for all  $j$ . Some of these roots may be irrational; some of these roots may be complex; and some of these roots may be repeated. Despite these complications, this theorem implies a unique representation of any polynomial of the form

$$p(x) = s \prod_{j=1}^n (x - r_j)$$

where the  $r_j$ 's are the roots and  $s$  is a scale factor. Once again, to represent a polynomial of degree  $n$ , we need a list of  $n + 1$  numbers: one scale factor and  $n$  roots.

Given a polynomial in this root representation, we can clearly evaluate it in  $O(n)$  time. Given two polynomials in root representation, we can easily multiply them in  $O(n)$  time by multiplying their scale factors and just concatenating the two root sequences.

Unfortunately, if we want to add two polynomials in root representation, we're out of luck. There's essentially *no* correlation between the roots of  $p$ , the roots of  $q$ , and the roots of  $p + q$ . We could

---

<sup>1</sup>I'm going to assume in this lecture that each arithmetic operation takes  $O(1)$  time. This may not be true in practice; in fact, one of the most powerful applications of FFTs is fast *integer* multiplication. One of the fastest integer multiplication algorithms, due to Schönhage and Strassen, multiplies two  $n$ -bit binary numbers using  $O(n \log n \log \log n \log \log \log n \dots 2^{O(\log^* n)})$  bit operations. The algorithm uses an  $n$ -element Fast Fourier Transform, which requires several  $O(\log n)$ -nit integer multiplications, which are carried out recursively (of course!), leading to the cascade of logs in the running time. The  $2^{O(\log^* n)}$  factor accounts for the constant factors accumulated at each of the  $O(\log^* n)$  levels of recursion. The *fastest* algorithm known for integer multiplication, published by Martin Fürer in 2007, uses  $O(n \log n 2^{O(\log^* n)})$  bit operations.

convert the polynomials to the more familiar coefficient representation first—this takes  $O(n^2)$  time using the high-school algorithms—but there's no easy way to convert the answer back. In fact, for most polynomials of degree 5 or more in coefficient form, it's *impossible* to compute roots exactly.<sup>2</sup>

### 2.2.2 Samples

Our third representation for polynomials comes from a different consequence of the Fundamental Theorem of Algebra. Given a list of  $n + 1$  pairs  $\{(x_0, y_0), (x_1, y_1), \dots, (x_n, y_n)\}$ , there is *exactly one* polynomial  $p$  of degree  $n$  such that  $p(x_j) = y_j$  for all  $j$ . This is just a generalization of the fact that any two points determine a unique line, because a line is the graph of a polynomial of degree 1. We say that the polynomial  $p$  **interpolates** the points  $(x_j, y_j)$ . As long as we agree on the sample locations  $x_j$  in advance, we once again need exactly  $n + 1$  numbers to represent a polynomial of degree  $n$ .

Adding or multiplying two polynomials in this sample representation is easy, as long as they use the same sample locations  $x_j$ . To add the polynomials, just add their sample values. To multiply two polynomials, just multiply their sample values; however, if we're multiplying two polynomials of degree  $n$ , we must *start* with  $2n + 1$  sample values for each polynomial, because that's how many we need to uniquely represent their product. Both algorithms run in  $O(n)$  time.

Unfortunately, evaluating a polynomial in this representation is no longer straightforward. The following formula, due to Lagrange, allows us to compute the value of any polynomial of degree  $n$  at any point, given a set of  $n + 1$  samples.

$$p(x) = \sum_{j=0}^{n-1} \left( \frac{y_j}{\prod_{k \neq j} (x_j - x_k)} \prod_{k \neq j} (x - x_k) \right)$$

Hopefully it's clear that formula actually describes a polynomial function of  $x$ , since each term in the rightmost sum is a scaled product of monomials. It's also not hard to verify that  $p(x_j) = y_j$  for every index  $j$ ; most of the terms of the sum vanish. As I mentioned earlier, the Fundamental Theorem of Algebra implies that  $p$  is *the only* polynomial that interpolates the points  $\{(x_j, y_j)\}$ . Lagrange's formula can be translated mechanically into an  $O(n^2)$ -time algorithm.

### 2.2.3 Summary

We find ourselves in the following frustrating situation. We have three representations for polynomials and three basic operations. Each representation allows us to almost trivially perform a different pair of operations in linear time, but the third takes at least quadratic time, if it can be done at all!

	evaluate	add	multiply
coefficients	$O(n)$	$O(n)$	<b><math>O(n^2)</math></b>
roots + scale	$O(n)$	$\infty$	$O(n)$
samples	<b><math>O(n^2)</math></b>	$O(n)$	$O(n)$

## 2.3 Converting Between Representations

What we need are fast algorithms to convert quickly from one representation to another. That way, when we need to perform an operation that's hard for our default representation, we can switch to a

---

<sup>2</sup>This is where numerical analysis comes from.

different representation that makes the operation easy, perform that operation, and then switch back. This strategy immediately rules out the root representation, since (as I mentioned earlier) finding roots of polynomials is impossible in general, at least if we're interested in exact results.

So how do we convert from coefficients to samples and back? Clearly, once we choose our sample positions  $x_j$ , we can compute each sample value  $y_j = p(x_j)$  in  $O(n)$  time from the coefficients using Horner's rule. So we can convert a polynomial of degree  $n$  from coefficients to samples in  $O(n^2)$  time. Lagrange's formula can be used to convert the sample representation back to the more familiar coefficient form. If we use the naïve algorithms for adding and multiplying polynomials (in coefficient form), this conversion takes  $O(n^3)$  time.

We can improve the cubic running time by observing that *both* conversion problems boil down to computing the product of a matrix and a vector. The explanation will be slightly simpler if we assume the polynomial has degree  $n - 1$ , so that  $n$  is the number of coefficients or samples. Fix a sequence  $x_0, x_1, \dots, x_{n-1}$  of sample positions, and let  $V$  be the  $n \times n$  matrix where  $v_{ij} = x_i^{j-1}$  (indexing rows and columns from 0 to  $n - 1$ ):

$$V = \begin{bmatrix} 1 & x_0 & x_0^2 & \cdots & x_0^{n-1} \\ 1 & x_1 & x_1^2 & \cdots & x_1^{n-1} \\ 1 & x_2 & x_2^2 & \cdots & x_2^{n-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n-1} & x_{n-1}^2 & \cdots & x_{n-1}^{n-1} \end{bmatrix}.$$

The matrix  $V$  is called a **Vandermonde matrix**. The vector of coefficients  $\vec{a} = (a_0, a_1, \dots, a_{n-1})$  and the vector of sample values  $\vec{y} = (y_0, y_1, \dots, y_{n-1})$  are related by the matrix equation

$$V\vec{a} = \vec{y},$$

or in more detail,

$$\begin{bmatrix} 1 & x_0 & x_0^2 & \cdots & x_0^{n-1} \\ 1 & x_1 & x_1^2 & \cdots & x_1^{n-1} \\ 1 & x_2 & x_2^2 & \cdots & x_2^{n-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n-1} & x_{n-1}^2 & \cdots & x_{n-1}^{n-1} \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ \vdots \\ a_{n-1} \end{bmatrix} = \begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ \vdots \\ y_{n-1} \end{bmatrix}.$$

Given this formulation, we can clearly transform any coefficient vector  $\vec{a}$  into the corresponding sample vector  $\vec{y}$  in  $O(n^2)$  time.

Conversely, if we know the sample values  $\vec{y}$ , we can recover the coefficients by solving a system of  $n$  linear equations in  $n$  unknowns; this takes  $O(n^3)$  time if we use Gaussian elimination.<sup>3</sup> But we can speed this up by implicitly hard-coding the sample positions into the algorithm. To convert from samples to coefficients, we can simply multiply the sample vector by the inverse of  $V$ , again in  $O(n^2)$  time.

$$\vec{a} = V^{-1}\vec{y}$$

Computing  $V^{-1}$  would take  $O(n^3)$  time if we had to do it from scratch using Gaussian elimination, but because we fixed the set of sample positions in advance, the matrix  $V^{-1}$  can be hard-coded directly into the algorithm.<sup>4</sup>

<sup>3</sup>In fact, Lagrange's formula is just a special case of Cramer's rule for solving linear systems.

<sup>4</sup>Actually, it is possible to invert an  $n \times n$  matrix in  $O(n^3)$  time, using fast matrix multiplication algorithms that closely resemble Karatsuba's sub-quadratic divide-and-conquer algorithm for integer/polynomial multiplication. On the other hand, my numerical-analysis colleagues have reasonable cause to shoot me for daring to suggest, even in passing, that anyone actually invert a matrix at all, ever.

So we can convert from coefficients to samples and back in  $O(n^2)$  time. At first glance, this result seems pointless; we can already add, multiply, or evaluate directly in either representation in  $O(n^2)$  time, so why bother? But there's a degree of freedom we haven't exploited—**We get to choose the sample positions!** Our conversion algorithm may be slow only because we're trying to be too general. If we choose a set of sample positions with the right recursive structure, we can perform this conversion more quickly.

## 2.4 Divide and Conquer

Any polynomial of degree  $n - 1$  can be expressed as a combination of two polynomials of degree  $(n/2) - 1$  as follows:

$$p(x) = p_{\text{even}}(x^2) + x \cdot p_{\text{odd}}(x^2).$$

The coefficients of  $p_{\text{even}}$  are just the even-degree coefficients of  $p$ , and the coefficients of  $p_{\text{odd}}$  are just the odd-degree coefficients of  $p$ . Thus, we can evaluate  $p(x)$  by recursively evaluating  $p_{\text{even}}(x^2)$  and  $p_{\text{odd}}(x^2)$  and performing  $O(1)$  additional arithmetic operations.

Now call a set  $X$  of  $n$  values **collapsing** if either of the following conditions holds:

- $X$  has one element.
- The set  $X^2 = \{x^2 \mid x \in X\}$  has exactly  $n/2$  elements and is (recursively) collapsing.

Clearly the size of any collapsing set is a power of 2. Given a polynomial  $p$  of degree  $n - 1$ , and a collapsing set  $X$  of size  $n$ , we can compute the set  $\{p(x) \mid x \in X\}$  of sample values as follows:

1. Recursively compute  $\{p_{\text{even}}(x^2) \mid x \in X\} = \{p_{\text{even}}(y) \mid y \in X^2\}$ .
2. Recursively compute  $\{p_{\text{odd}}(x^2) \mid x \in X\} = \{p_{\text{odd}}(y) \mid y \in X^2\}$ .
3. For each  $x \in X$ , compute  $p(x) = p_{\text{even}}(x^2) + x \cdot p_{\text{odd}}(x^2)$ .

The running time of this algorithm satisfies the familiar recurrence  $T(n) = 2T(n/2) + \Theta(n)$ , which as we all know solves to  $T(n) = \Theta(n \log n)$ .

Great! Now all we need is a sequence of arbitrarily large collapsible sets. The simplest method to construct such sets is just to invert the recursive definition: If  $X$  is a collapsible set of size  $n$  that does not contain the number 0, then  $\sqrt{X} = \{\pm\sqrt{x} \mid x \in X\}$  is a collapsible set of size  $2n$ . This observation gives us an infinite sequence of collapsible sets, starting as follows:<sup>5</sup>

$$\begin{aligned} X_1 &:= \{1\} \\ X_2 &:= \{1, -1\} \\ X_4 &:= \{1, -1, i, -i\} \\ X_8 &:= \left\{1, -1, i, -i, \frac{\sqrt{2}}{2} + \frac{\sqrt{2}}{2}i, -\frac{\sqrt{2}}{2} - \frac{\sqrt{2}}{2}i, \frac{\sqrt{2}}{2} - \frac{\sqrt{2}}{2}i, -\frac{\sqrt{2}}{2} + \frac{\sqrt{2}}{2}i\right\} \end{aligned}$$

---

<sup>5</sup>In this lecture,  $i$  always represents the square root of  $-1$ . Computer scientists are used to thinking of  $i$  as an integer index into a sequence, an array, or a for-loop, but we obviously can't do that here. The physicist's habit of using  $j = \sqrt{-1}$  just delays the problem (How do physicists write quaternions?), and typographical tricks like  $I$  or  $\mathbf{i}$  or Mathematica's  $\mathbf{i}$  are just stupid.

## 2.5 The Discrete Fourier Transform

For any  $n$ , the elements of  $X_n$  are called the **complex  $n$ th roots of unity**; these are the roots of the polynomial  $x^n - 1 = 0$ . These  $n$  complex values are spaced exactly evenly around the unit circle in the complex plane. Every  $n$ th root of unity is a power of the primitive  $n$ th root

$$\omega_n = e^{2\pi i/n} = \cos \frac{2\pi}{n} + i \sin \frac{2\pi}{n}.$$

A typical  $n$ th root of unity has the form

$$\omega_n^k = e^{(2\pi i/n)k} = \cos \left( \frac{2\pi}{n} k \right) + i \sin \left( \frac{2\pi}{n} k \right).$$

These complex numbers have several useful properties for any integers  $n$  and  $k$ :

- There are exactly  $n$  different  $n$ th roots of unity:  $\omega_n^k = \omega_n^{k \bmod n}$ .
- If  $n$  is even, then  $\omega_n^{k+n/2} = -\omega_n^k$ ; in particular,  $\omega_n^{n/2} = -\omega_n^0 = -1$ .
- $1/\omega_n^k = \omega_n^{-k} = \overline{\omega_n^k} = (\overline{\omega_n})^k$ , where the bar represents complex conjugation:  $\overline{a+bi} = a-bi$
- $\omega_n = \omega_{kn}^k$ . Thus, every  $n$ th root of unity is also a  $(kn)$ th root of unity.

These properties imply immediately that if  $n$  is a power of 2, then the set of all  $n$ th roots of unity is collapsible!

If we sample a polynomial of degree  $n-1$  at the  $n$ th roots of unity, the resulting list of sample values is called the **discrete Fourier transform** of the polynomial (or more formally, of its coefficient vector). Thus, given an array  $P[0..n-1]$  of coefficients, its discrete Fourier transform is the vector  $P^*[0..n-1]$  defined as follows:

$$P^*[j] := p(\omega_n^j) = \sum_{k=0}^{n-1} P[k] \cdot \omega_n^{jk}$$

As we already observed, the fact that sets of roots of unity are collapsible implies that we can compute the discrete Fourier transform in  $O(n \log n)$  time. The resulting algorithm, called the **fast Fourier transform**, was popularized by Cooley and Tukey in 1965.<sup>6</sup> The algorithm assumes that  $n$  is a power of two; if necessary, we can just pad the coefficient vector with zeros.

---

<sup>6</sup>Actually, Cooley and Tukey's FFT algorithm was previously published by Runge and König in 1924, and again by Yates in 1932, and again by Stumpf in 1937, and again by Danielson and Lanczos in 1942. So of course it's often called the **Cooley-Tukey algorithm**. But the algorithm was first used by Gauss in the 1800s for calculating the paths of asteroids from a finite number of equally-spaced observations. By hand. Fourier himself always did it the hard way.

Tukey apparently developed the algorithm to help detect Soviet nuclear tests without actually visiting Soviet nuclear facilities, by interpolating off-shore seismic readings. Without their rediscovery of the FFT algorithm, the nuclear test ban treaty would never have been ratified, and we'd all be speaking Russian, or more likely, whatever language radioactive glass speaks.

```

 $\text{FFT}(P[0..n-1]):$ 
  if  $n = 1$ 
    return  $P$ 
  for  $j \leftarrow 0$  to  $n/2 - 1$ 
     $U[j] \leftarrow P[2j]$ 
     $V[j] \leftarrow P[2j + 1]$ 
     $U^* \leftarrow \text{FFT}(U[0..n/2-1])$ 
     $V^* \leftarrow \text{FFT}(V[0..n/2-1])$ 
     $\omega_n \leftarrow \cos(\frac{2\pi}{n}) + i \sin(\frac{2\pi}{n})$ 
     $\omega \leftarrow 1$ 
    for  $j \leftarrow 0$  to  $n/2 - 1$ 
       $P^*[j] \leftarrow U^*[j] + \omega \cdot V^*[j]$ 
       $P^*[j+n/2] \leftarrow U^*[j] - \omega \cdot V^*[j]$ 
       $\omega \leftarrow \omega \cdot \omega_n$ 
  return  $P^*[0..n-1]$ 

```

## 2.6 Inverting the FFT

We also need to recover the coefficients of the product from the new sample values. Recall that the transformation from coefficients to sample values is *linear*; the sample vector is the product of a Vandermonde matrix  $V$  and the coefficient vector. For the discrete Fourier transform, each entry in  $V$  is an  $n$ th root of unity; specifically,

$$v_{jk} = \omega_n^{jk}$$

for all integers  $j$  and  $k$ . Thus,

$$V = \begin{bmatrix} 1 & 1 & 1 & 1 & \dots & 1 \\ 1 & \omega_n & \omega_n^2 & \omega_n^3 & \dots & \omega_n^{n-1} \\ 1 & \omega_n^2 & \omega_n^4 & \omega_n^6 & \dots & \omega_n^{2(n-1)} \\ 1 & \omega_n^3 & \omega_n^6 & \omega_n^9 & \dots & \omega_n^{3(n-1)} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \omega_n^{n-1} & \omega_n^{2(n-1)} & \omega_n^{3(n-1)} & \dots & \omega_n^{(n-1)^2} \end{bmatrix}$$

To invert the discrete Fourier transform, converting sample values back to coefficients, we just have to multiply  $P^*$  by the inverse matrix  $V^{-1}$ . The following amazing fact implies that this is almost the same as multiplying by  $V$  itself:

**Claim:**  $V^{-1} = \bar{V}/n$

**Proof:** Let  $W = \bar{V}/n$ . We just have to show that  $M = VW$  is the identity matrix. We can compute a single entry in  $M$  as follows:

$$m_{jk} = \sum_{l=0}^{n-1} v_{jl} \cdot w_{lk} = \sum_{l=0}^{n-1} \omega_n^{jl} \cdot \bar{\omega}_n^{-lk} / n = \frac{1}{n} \sum_{l=0}^{n-1} \omega_n^{jl-lk} = \frac{1}{n} \sum_{l=0}^{n-1} (\omega_n^{j-k})^l$$

If  $j = k$ , then  $\omega_n^{j-k} = \omega_n^0 = 1$ , so

$$m_{jk} = \frac{1}{n} \sum_{l=0}^{n-1} 1 = \frac{n}{n} = 1,$$

and if  $j \neq k$ , we have a geometric series

$$m_{jk} = \sum_{l=0}^{n-1} (\omega_n^{j-k})^l = \frac{(\omega_n^{j-k})^n - 1}{\omega_n^{j-k} - 1} = \frac{(\omega_n^n)^{j-k} - 1}{\omega_n^{j-k} - 1} = \frac{1^{j-k} - 1}{\omega_n^{j-k} - 1} = 0.$$

That's it!  $\square$

In other words, if  $W = V^{-1}$  then  $w_{jk} = \overline{v_{jk}}/n = \overline{\omega_n^{jk}}/n = \omega_n^{-jk}/n$ . What this means for us computer scientists is that any algorithm for computing the discrete Fourier transform can be easily modified to compute the inverse transform as well.

```
INVERSEFFT( $P^*[0..n-1]$ ):
  if  $n = 1$ 
    return  $P$ 
  for  $j \leftarrow 0$  to  $n/2 - 1$ 
     $U^*[j] \leftarrow P^*[2j]$ 
     $V^*[j] \leftarrow P^*[2j + 1]$ 
     $U \leftarrow \text{INVERSEFFT}(U[0..n/2 - 1])$ 
     $V \leftarrow \text{INVERSEFFT}(V[0..n/2 - 1])$ 
     $\overline{\omega_n} \leftarrow \cos(\frac{2\pi}{n}) - i \sin(\frac{2\pi}{n})$ 
     $\omega \leftarrow 1$ 
    for  $j \leftarrow 0$  to  $n/2 - 1$ 
       $P[j] \leftarrow 2(U[j] + \omega \cdot V[j])$ 
       $P[j + n/2] \leftarrow 2(U[j] - \omega \cdot V[j])$ 
       $\omega \leftarrow \omega \cdot \overline{\omega_n}$ 
  return  $P[0..n-1]$ 
```

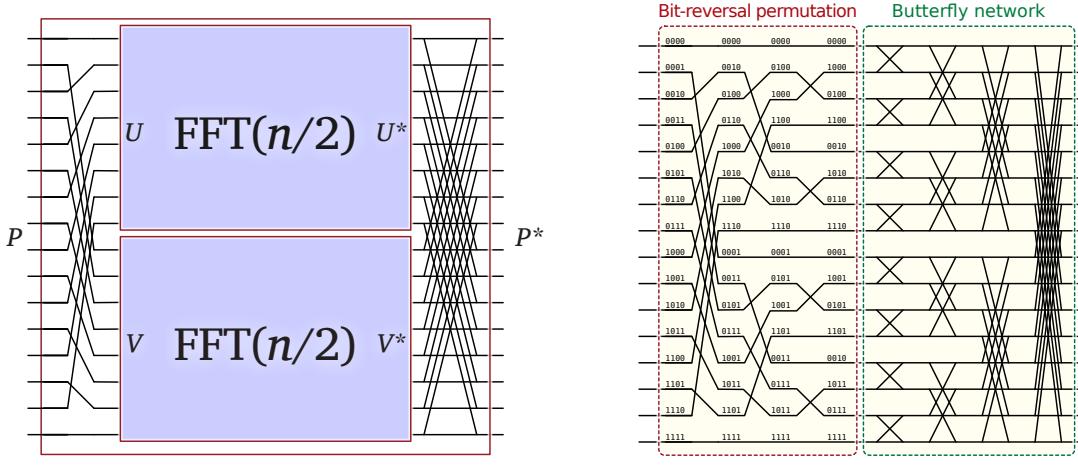
## 2.7 Fast Multiplication

Finally, given two polynomials  $p$  and  $q$ , each represented by an array of coefficients, we can multiply them in  $\Theta(n \log n)$  arithmetic operations as follows. First, pad the coefficient vectors and with zeros until the size is a power of two greater than or equal to the sum of the degrees. Then compute the DFTs of each coefficient vector, multiply the sample values one by one, and compute the inverse DFT of the resulting sample vector.

```
FFTMULTIPLY( $P[0..n-1], Q[0..m-1]$ ):
   $\ell \leftarrow \lceil \lg(n+m) \rceil$ 
  for  $j \leftarrow n$  to  $2^\ell - 1$ 
     $P[j] \leftarrow 0$ 
  for  $j \leftarrow m$  to  $2^\ell - 1$ 
     $Q[j] \leftarrow 0$ 
   $P^* \leftarrow FFT(P)$ 
   $Q^* \leftarrow FFT(Q)$ 
  for  $j \leftarrow 0$  to  $2^\ell - 1$ 
     $R^*[j] \leftarrow P^*[j] \cdot Q^*[j]$ 
  return INVERSEFFT( $R^*$ )
```

## 2.8 Inside the FFT

FFTs are often implemented in hardware as circuits. To see the recursive structure of the circuit, let's connect the top-level inputs and outputs to the inputs and outputs of the recursive calls. On the left we split the input  $P$  into two recursive inputs  $U$  and  $V$ . On the right, we combine the outputs  $U^*$  and  $V^*$  to obtain the final output  $P^*$ .



The recursive structure of the FFT algorithm.

If we expand this recursive structure completely, we see that the circuit splits naturally into two parts. The left half computes the *bit-reversal permutation* of the input. To find the position of  $P[k]$  in this permutation, write  $k$  in binary, and then read the bits backward. For example, in an 8-element bit-reversal permutation,  $P[3] = P[011_2]$  ends up in position 6 =  $110_2$ . The right half of the FFT circuit is a *butterfly network*. Butterfly networks are often used to route between processors in massively-parallel computers, because they allow any two processors to communicate in only  $O(\log n)$  steps.

## Exercises

1. For any two sets  $X$  and  $Y$  of integers, the Minkowski sum  $X + Y$  is the set of all pairwise sums  $\{x + y \mid x \in X, y \in Y\}$ .
  - (a) Describe an analyze and algorithm to compute the number of elements in  $X + Y$  in  $O(n^2 \log n)$  time. [Hint: The answer is **not** always  $n^2$ .]
  - (b) Describe and analyze an algorithm to compute the number of elements in  $X + Y$  in  $O(M \log M)$  time, where  $M$  is the largest absolute value of any element of  $X \cup Y$ . [Hint: What's this lecture about?]
2. (a) Describe an algorithm that determines whether a given set of  $n$  integers contains two elements whose sum is zero, in  $O(n \log n)$  time.
  - (b) Describe an algorithm that determines whether a given set of  $n$  integers contains *three* elements whose sum is zero, in  $O(n^2)$  time.
  - (c) Now suppose the input set  $X$  contains only integers between  $-10000n$  and  $10000n$ . Describe an algorithm that determines whether  $X$  contains three elements whose sum is zero, in  $O(n \log n)$  time. [Hint: Hint.]

'Tis a lesson you should heed,  
 Try, try again;  
 If at first you don't succeed,  
 Try, try again;  
 Then your courage should appear,  
 For, if you will persevere,  
 You will conquer, never fear;  
 Try, try again.

— Thomas H. Palmer, *The Teacher's Manual: Being an Exposition of an Efficient and Economical System of Education Suited to the Wants of a Free People* (1840)

*When you come to a fork in the road, take it.*

— Yogi Berra

## 3 Backtracking

In this lecture, I want to describe another recursive algorithm strategy called **backtracking**. A backtracking algorithm tries to build a solution to a computational problem incrementally. Whenever the algorithm needs to decide between multiple alternatives to the next component of the solution, it simply tries all possible options recursively.

### 3.1 *n* Queens

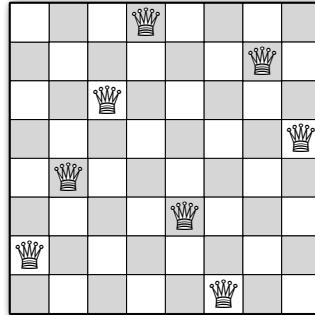
The prototypical backtracking problem is the classical ***n* Queens Problem**, first proposed by German chess enthusiast Max Bezzel in 1848 for the standard  $8 \times 8$  board, and both solved and generalized to larger boards by Franz Nauck in 1850. The problem is to place  $n$  queens on an  $n \times n$  chessboard, so that no two queens can attack each other. For readers not familiar with the rules of chess, this means that no two queens are in the same row, column, or diagonal.

Obviously, in any solution to the  $n$ -Queens problem, there is exactly one queen in each row. So we will represent our possible solutions using an array  $Q[1..n]$ , where  $Q[i]$  indicates which square in row  $i$  contains a queen, or 0 if no queen has yet been placed in row  $i$ . To find a solution, we put queens on the board row by row, starting at the top. A *partial* solution is an array  $Q[1..n]$  whose first  $r - 1$  entries are positive and whose last  $n - r + 1$  entries are all zeros, for some integer  $r$ .

The following recursive algorithm recursively enumerates all complete  $n$ -queens solutions that are consistent with a given partial solution. The input parameter  $r$  is the first empty row. Thus, to compute all  $n$ -queens solutions with no restrictions, we would call `RECURSIVENQUEENS(Q[1..n], 1)`.

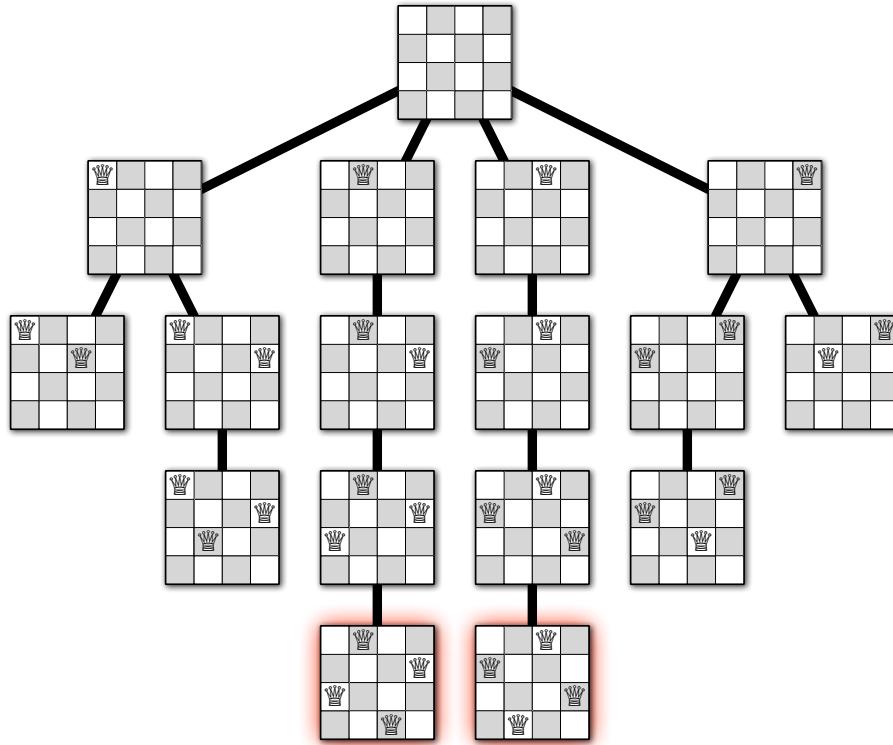
```

RECURSIVENQUEENS(Q[1..n], r):
  if  $r = n + 1$ 
    print Q
  else
    for  $j \leftarrow 1$  to  $n$ 
      legal  $\leftarrow$  TRUE
      for  $i \leftarrow 1$  to  $r - 1$ 
        if  $(Q[i] = j) \text{ or } (Q[i] = j + r - i) \text{ or } (Q[i] = j - r + i)$ 
          legal  $\leftarrow$  FALSE
      if legal
         $Q[r] \leftarrow j$ 
        RECURSIVENQUEENS(Q[1..n],  $r + 1$ )
    
```



One solution to the 8 queens problem, represented by the array [4,7,3,8,2,5,1,6]

Like most recursive algorithms, the execution of a backtracking algorithm can be illustrated using a **recursion tree**. The root of the recursion tree corresponds to the original invocation of the algorithm; edges in the tree correspond to recursive calls. A path from the root down to any node shows the history of a partial solution to the  $n$ -Queens problem, as queens are added to successive rows. The leaves correspond to partial solutions that cannot be extended, either because there is already a queen on every row, or because every position in the next empty row is in the same row, column, or diagonal as an existing queen. The backtracking algorithm simply performs a traversal of this tree.



The complete recursion tree for our algorithm for the 4 queens problem.

### 3.2 Subset Sum

Let's consider a more complicated problem, called **SUBSETSUM**: Given a set  $X$  of positive integers and target integer  $T$ , is there a subset of elements in  $X$  that add up to  $T$ ? Notice that there can be more than one such subset. For example, if  $X = \{8, 6, 7, 5, 3, 10, 9\}$  and  $T = 15$ , the answer is TRUE, thanks to

the subsets  $\{8, 7\}$  or  $\{7, 5, 3\}$  or  $\{6, 9\}$  or  $\{5, 10\}$ . On the other hand, if  $X = \{11, 6, 5, 1, 7, 13, 12\}$  and  $T = 15$ , the answer is FALSE.

There are two trivial cases. If the target value  $T$  is zero, then we can immediately return TRUE, because empty set is a subset of *every* set  $X$ , and the elements of the empty set add up to zero.<sup>1</sup> On the other hand, if  $T < 0$ , or if  $T \neq 0$  but the set  $X$  is empty, then we can immediately return FALSE.

For the general case, consider an arbitrary element  $x \in X$ . (We've already handled the case where  $X$  is empty.) There is a subset of  $X$  that sums to  $T$  if and only if one of the following statements is true:

- There is a subset of  $X$  that *includes*  $x$  and whose sum is  $T$ .
- There is a subset of  $X$  that *excludes*  $x$  and whose sum is  $T$ .

In the first case, there must be a subset of  $X \setminus \{x\}$  that sums to  $T - x$ ; in the second case, there must be a subset of  $X \setminus \{x\}$  that sums to  $T$ . So we can solve  $\text{SUBSETSUM}(X, T)$  by reducing it to two simpler instances:  $\text{SUBSETSUM}(X \setminus \{x\}, T - x)$  and  $\text{SUBSETSUM}(X \setminus \{x\}, T)$ . Here's how the resulting recursive algorithm might look if  $X$  is stored in an array.

```
SUBSETSUM( $X[1..n]$ ,  $T$ ):
  if  $T = 0$ 
    return TRUE
  else if  $T < 0$  or  $n = 0$ 
    return FALSE
  else
    return ( $\text{SUBSETSUM}(X[2..n], T) \vee \text{SUBSETSUM}(X[2..n], T - X[1])$ )
```

Proving this algorithm correct is a straightforward exercise in induction. If  $T = 0$ , then the elements of the empty subset sum to  $T$ , so TRUE is the correct output. Otherwise, if  $T$  is negative or the set  $X$  is empty, then no subset of  $X$  sums to  $T$ , so FALSE is the correct output. Otherwise, if there is a subset that sums to  $T$ , then either it contains  $X[1]$  or it doesn't, and the Recursion Fairy correctly checks for each of those possibilities. Done.

The running time  $T(n)$  clearly satisfies the recurrence  $T(n) \leq 2T(n - 1) + O(1)$ , which we can solve using either recursion trees or annihilators (or just guessing) to obtain the upper bound  $T(n) = O(2^n)$ . In the worst case, the recursion tree for this algorithm is a complete binary tree with depth  $n$ .

Here is a similar recursive algorithm that actually *constructs* a subset of  $X$  that sums to  $T$ , if one exists. This algorithm also runs in  $O(2^n)$  time.

```
CONSTRUCTSUBSET( $X[1..n]$ ,  $T$ ):
  if  $T = 0$ 
    return  $\emptyset$ 
  if  $T < 0$  or  $n = 0$ 
    return NONE
   $Y \leftarrow \text{CONSTRUCTSUBSET}(X[2..n], T)$ 
  if  $Y \neq \text{NONE}$ 
    return  $Y$ 
   $Y \leftarrow \text{CONSTRUCTSUBSET}(X[2..n], T - X[1])$ 
  if  $Y \neq \text{NONE}$ 
    return  $Y \cup \{X[1]\}$ 
  return NONE
```

---

<sup>1</sup>There's no base case like the vacuous base case!

### 3.3 Longest Increasing Subsequence

Now suppose we are given a sequence of integers, and we want to find the longest subsequence whose elements are in increasing order. More concretely, the input is an array  $A[1..n]$  of integers, and we want to find the longest sequence of indices  $1 \leq i_1 < i_2 < \dots < i_k \leq n$  such that  $A[i_j] < A[i_{j+1}]$  for all  $j$ .

To derive a recursive algorithm for this problem, we start with a recursive definition of the kinds of objects we're playing with: sequences and subsequences.

A sequence of integers is either empty  
or an integer followed by a sequence of integers.

This definition suggests the following strategy for devising a recursive algorithm. If the input sequence is empty, there's nothing to do. Otherwise, we only need to figure out what to do with the first element of the input sequence; the Recursion Fairy will take care of everything else. We can formalize this strategy somewhat by giving a recursive definition of subsequence (using array notation to represent sequences):

The only subsequence of the empty sequence is the empty sequence.  
A subsequence of  $A[1..n]$  is either a subsequence of  $A[2..n]$   
or  $A[1]$  followed by a subsequence of  $A[2..n]$ .

We're not just looking for just *any* subsequence, but a *longest* subsequence with the property that elements are in *increasing* order. So let's try to add those two conditions to our definition. (I'll omit the familiar vacuous base case.)

The LIS of  $A[1..n]$  is  
either the LIS of  $A[2..n]$   
or  $A[1]$  followed by the LIS of  $A[2..n]$  with elements larger than  $A[1]$ ,  
whichever is longer.

This definition is correct, but it's not quite recursive—we're defining the object 'longest increasing subsequence' in terms of the slightly *different* object 'longest increasing subsequence with elements larger than  $x$ ', which we haven't properly defined yet. Fortunately, this second object has a very similar recursive definition. (Again, I'm omitting the vacuous base case.)

If  $A[1] \leq x$ , the LIS of  $A[1..n]$  with elements larger than  $x$  is  
the LIS of  $A[2..n]$  with elements larger than  $x$ .  
Otherwise, the LIS of  $A[1..n]$  with elements larger than  $x$  is  
either the LIS of  $A[2..n]$  with elements larger than  $x$   
or  $A[1]$  followed by the LIS of  $A[2..n]$  with elements larger than  $A[1]$ ,  
whichever is longer.

The longest increasing subsequence without restrictions can now be redefined as the longest increasing subsequence with elements larger than  $-\infty$ . Rewriting this recursive definition into pseudocode gives us the following recursive algorithm.

```
LIS(A[1..n]):  
    return LISBIGGER(-∞, A[1..n])
```

```
LISBIGGER(prev, A[1..n]):  
    if n = 0  
        return 0  
    else  
        max ← LISBIGGER(prev, A[2..n])  
        if A[1] > prev  
            L ← 1 + LISBIGGER(A[1], A[2..n])  
            if L > max  
                max ← L  
        return max
```

The running time of this algorithm satisfies the recurrence  $T(n) \leq 2T(n-1) + O(1)$ , which as usual implies that  $T(n) = O(2^n)$ . We really shouldn't be surprised by this running time; in the worst case, the algorithm examines each of the  $2^n$  subsequences of the input array.

The following alternative strategy avoids defining a new object with the ‘larger than  $x$ ’ constraint. We still only have to decide whether to include or exclude the first element  $A[1]$ . We consider the case where  $A[1]$  is excluded exactly the same way, but to consider the case where  $A[1]$  is included, we remove any elements of  $A[2..n]$  that are larger than  $A[1]$  before we recurse. This new strategy gives us the following algorithm:

```
FILTER(A[1..n], x):  
    j ← 1  
    for i ← 1 to n  
        if A[i] > x  
            B[j] ← A[i]; j ← j + 1  
    return B[1..j]
```

```
LIS(A[1..n]):  
    if n = 0  
        return 0  
    else  
        max ← LIS(prev, A[2..n])  
        L ← 1 + LIS(A[1], FILTER(A[2..n], A[1]))  
        if L > max  
            max ← L  
    return max
```

The FILTER subroutine clearly runs in  $O(n)$  time, so the running time of LIS satisfies the recurrence  $T(n) \leq 2T(n-1) + O(n)$ , which solves to  $T(n) \leq O(2^n)$  by the annihilator method. This upper bound pessimistically assumes that FILTER never actually removes any elements; indeed, if the input sequence is sorted in increasing order, this assumption is correct.

### 3.4 Optimal Binary Search Trees

Our next example combines recursive backtracking with the divide-and-conquer strategy.

Recall that the running time for a successful search in a binary search tree is proportional to the number of ancestors of the target node.<sup>2</sup> As a result, the worst-case search time is proportional to the depth of the tree. Thus, to minimize the worst-case search time, the height of the tree should be as small as possible; by this metric, the ideal tree is perfectly balanced.

In many applications of binary search trees, however, it is more important to minimize the total cost of several searches rather than the worst-case cost of a single search. If  $x$  is a more ‘popular’ search target than  $y$ , we can save time by building a tree where the depth of  $x$  is smaller than the depth of  $y$ , even if that means increasing the overall depth of the tree. A perfectly balanced tree is *not* the best choice if some items are significantly more popular than others. In fact, a totally unbalanced tree of depth  $\Omega(n)$  might actually be the best choice!

<sup>2</sup>An *ancestor* of a node  $v$  is either the node itself or an ancestor of the parent of  $v$ . A *proper* ancestor of  $v$  is either the parent of  $v$  or a proper ancestor of the parent of  $v$ .

This situation suggests the following problem. Suppose we are given a sorted array of  $n$  keys  $A[1..n]$  and an array of corresponding *access frequencies*  $f[1..n]$ . Our task is to build the binary search tree that minimizes the *total* search time, assuming that there will be exactly  $f[i]$  searches for each key  $A[i]$ .

Before we think about how to solve this problem, we should first come up with a good recursive definition of the function we are trying to optimize! Suppose we are also given a binary search tree  $T$  with  $n$  nodes. Let  $v_i$  denote the node that stores  $A[i]$ , and let  $r$  be the index of the root node. Ignoring constant factors, the cost of searching for  $A[i]$  is the number of nodes on the path from the root  $v_r$  to  $v_i$ . Thus, the total cost of performing all the binary searches is given by the following expression:

$$\text{Cost}(T, f[1..n]) = \sum_{i=1}^n f[i] \cdot \#\text{nodes between } v_r \text{ and } v_i$$

Every search path includes the root node  $v_r$ . If  $i < r$ , then all other nodes on the search path to  $v_i$  are in the left subtree; similarly, if  $i > r$ , all other nodes on the search path to  $v_i$  are in the right subtree. Thus, we can partition the cost function into three parts as follows:

$$\begin{aligned} \text{Cost}(T, f[1..n]) &= \sum_{i=1}^{r-1} f[i] \cdot \#\text{nodes between } \text{left}(v_r) \text{ and } v_i \\ &\quad + \sum_{i=1}^n f[i] \\ &\quad + \sum_{i=r+1}^n f[i] \cdot \#\text{nodes between } \text{right}(v_r) \text{ and } v_i \end{aligned}$$

Now the first and third summations look exactly like our original expression (\*) for  $\text{Cost}(T, f[1..n])$ . Simple substitution gives us our recursive definition for  $\text{Cost}$ :

$$\text{Cost}(T, f[1..n]) = \text{Cost}(\text{left}(T), f[1..r-1]) + \sum_{i=1}^n f[i] + \text{Cost}(\text{right}(T), f[r+1..n])$$

The base case for this recurrence is, as usual,  $n = 0$ ; the cost of performing no searches in the empty tree is zero.

Now our task is to compute the tree  $T_{\text{opt}}$  that minimizes this cost function. Suppose we somehow magically knew that the root of  $T_{\text{opt}}$  is  $v_r$ . Then the recursive definition of  $\text{Cost}(T, f)$  immediately implies that the left subtree  $\text{left}(T_{\text{opt}})$  must be the optimal search tree for the keys  $A[1..r-1]$  and access frequencies  $f[1..r-1]$ . Similarly, the right subtree  $\text{right}(T_{\text{opt}})$  must be the optimal search tree for the keys  $A[r+1..n]$  and access frequencies  $f[r+1..n]$ . **Once we choose the correct key to store at the root, the Recursion Fairy will automatically construct the rest of the optimal tree for us.** More formally, let  $\text{OptCost}(f[1..n])$  denote the total cost of the optimal search tree for the given frequency counts. We immediately have the following recursive definition.

$$\text{OptCost}(f[1..n]) = \min_{1 \leq r \leq n} \left\{ \text{OptCost}(f[1..r-1]) + \sum_{i=1}^n f[i] + \text{OptCost}(f[r+1..n]) \right\}$$

Again, the base case is  $\text{OptCost}(f[1..0]) = 0$ ; the best way to organize no keys, which we will plan to search zero times, is by storing them in the empty tree!

This recursive definition can be translated mechanically into a recursive algorithm, whose running time  $T(n)$  satisfies the recurrence

$$T(n) = \Theta(n) + \sum_{k=1}^n (T(k-1) + T(n-k)).$$

The  $\Theta(n)$  term comes from computing the total number of searches  $\sum_{i=1}^n f[i]$ .

Yeah, that's one ugly recurrence, but it's actually easier to solve than it looks. To transform it into a more familiar form, we regroup and collect identical terms, subtract the recurrence for  $T(n-1)$  to get rid of the summation, and then regroup again.

$$\begin{aligned} T(n) &= \Theta(n) + 2 \sum_{k=0}^{n-1} T(k) \\ T(n-1) &= \Theta(n-1) + 2 \sum_{k=0}^{n-2} T(k) \\ T(n) - T(n-1) &= \Theta(1) + 2T(n-1) \\ T(n) &= 3T(n-1) + \Theta(1) \end{aligned}$$

The solution  $T(n) = \Theta(3^n)$  now follows from the annihilator method.

Let me emphasize that this recursive algorithm does *not* examine all possible binary search trees. The number of binary search trees with  $n$  nodes satisfies the recurrence

$$N(n) = \sum_{r=1}^{n-1} (N(r-1) \cdot N(n-r)),$$

which has the closed-form solution  $N(n) = \Theta(4^n / \sqrt{n})$ . Our algorithm saves considerable time by searching *independently* for the optimal left and right subtrees. A full enumeration of binary search trees would consider all possible *pairings* of left and right subtrees; hence the product in the recurrence for  $N(n)$ .

## Exercises

1. (a) Let  $A[1..m]$  and  $B[1..n]$  be two arbitrary arrays. A *common subsequence* of  $A$  and  $B$  is both a subsequence of  $A$  and a subsequence of  $B$ . Give a simple recursive definition for the function  $lcs(A, B)$ , which gives the length of the *longest* common subsequence of  $A$  and  $B$ .
- (b) Let  $A[1..m]$  and  $B[1..n]$  be two arbitrary arrays. A *common supersequence* of  $A$  and  $B$  is another sequence that contains both  $A$  and  $B$  as subsequences. Give a simple recursive definition for the function  $scs(A, B)$ , which gives the length of the *shortest* common supersequence of  $A$  and  $B$ .
- (c) Call a sequence  $X[1..n]$  *oscillating* if  $X[i] < X[i+1]$  for all even  $i$ , and  $X[i] > X[i+1]$  for all odd  $i$ . Give a simple recursive definition for the function  $los(A)$ , which gives the length of the longest oscillating subsequence of an arbitrary array  $A$  of integers.
- (d) Give a simple recursive definition for the function  $sos(A)$ , which gives the length of the shortest oscillating supersequence of an arbitrary array  $A$  of integers.

- (e) Call a sequence  $X[1..n]$  *accelerating* if  $2 \cdot X[i] < X[i - 1] + X[i + 1]$  for all  $i$ . Give a simple recursive definition for the function  $lxs(A)$ , which gives the length of the longest accelerating subsequence of an arbitrary array  $A$  of integers.

**For more backtracking exercises, see the next two lecture notes!**

*Wouldn't the sentence "I want to put a hyphen between the words Fish and And and And and Chips in my Fish-And-Chips sign." have been clearer if quotation marks had been placed before Fish, and between Fish and and, and and and Chips, as well as after Chips?*<sup>1</sup>

— Martin Gardner, *Aha! Insight* (1978)

## \*4 Efficient Exponential-Time Algorithms

In another lecture note, we discuss the class of **NP-hard** problems. For every problem in this class, the fastest algorithm anyone knows has an exponential running time. Moreover, there is very strong evidence (but alas, no proof) that it is *impossible* to solve any NP-hard problem in less than exponential time—it's not that we're all stupid; the problems really are that hard! Unfortunately, an enormous number of problems that arise in practice are NP-hard; for some of these problems, even *approximating* the right answer is NP-hard.

Suppose we absolutely have to find the exact solution to some NP-hard problem. A polynomial-time algorithm is almost certainly out of the question; the best running time we can hope for is exponential. But *which* exponential? An algorithm that runs in  $O(1.5^n)$  time, while still unusable for large problems, is still significantly better than an algorithm that runs in  $O(2^n)$  time!

For most NP-hard problems, the only approach that is guaranteed to find an optimal solution is recursive backtracking. The most straightforward version of this approach is to recursively generate *all* possible solutions and check each one: all satisfying assignments, or all vertex colorings, or all subsets, or all permutations, or whatever. However, most NP-hard problems have some additional structure that allows us to prune away most of the branches of the recursion tree, thereby drastically reducing the running time.

### 4.1 3SAT

Let's consider the mother of all NP-hard problems: 3SAT. Given a boolean formula in conjunctive normal form, with at most three literals in each clause, our task is to determine whether any assignment of values of the variables makes the formula true. Yes, this problem is NP-hard, which means that a polynomial algorithm is almost certainly impossible. Too bad; we have to solve the problem anyway.

The trivial solution is to try every possible assignment. We'll evaluate the running time of our 3SAT algorithms in terms of the number of variables in the formula, so let's call that  $n$ . Provided any clause appears in our input formula at most once—a condition that we can easily enforce in polynomial time—the overall input size is  $O(n^3)$ . There are  $2^n$  possible assignments, and we can evaluate each assignment in  $O(n^3)$  time, so the overall running time is  $O(2^n n^3)$ .

Since polynomial factors like  $n^3$  are essentially noise when the overall running time is exponential, from now on I'll use  $\text{poly}(n)$  to represent some arbitrary polynomial in  $n$ ; in other words,  $\text{poly}(n) = n^{O(1)}$ . For example, the trivial algorithm for 3SAT runs in time  $O(2^n \text{poly}(n))$ .

---

<sup>1</sup>If you ever decide to read this sentence out loud, be sure to pause briefly between 'Fish and and' and 'and and and And', 'and and and And' and 'and And and and', 'and And and and' and 'and and and And', 'and and and And' and 'and And and and', and 'and And and and' and 'and and and Chips'!

Did you notice the punctuation I carefully inserted between 'Fish and and' and 'and', 'and' and 'and and and And', 'and and and And' and 'and and And', 'and and And' and 'and', 'and' and 'and And and and', 'and And and and' and 'and And and and', 'and And and and' and 'and', 'and' and 'and and and And', 'and and and And' and 'and and and And', 'and and and And' and 'and', 'and' and 'and And and and', 'and And and and' and 'and', 'and' and 'and And and and', 'and And and and' and 'and', and 'and' and 'and And and and', and 'and', and 'and' and 'and and and Chips'?

We can make this algorithm smarter by exploiting the special recursive structure of 3CNF formulas:

A 3CNF formula is either nothing  
or a clause with three literals  $\wedge$  a 3CNF formula

Suppose we want to decide whether some 3CNF formula  $\Phi$  with  $n$  variables is satisfiable. Of course this is trivial if  $\Phi$  is the empty formula, so suppose

$$\Phi = (x \vee y \vee z) \wedge \Phi'$$

for some literals  $x, y, z$  and some 3CNF formula  $\Phi'$ . By distributing the  $\wedge$  across the  $\vee$ s, we can rewrite  $\Phi$  as follows:

$$\Phi = (x \wedge \Phi') \vee (y \wedge \Phi') \vee (z \wedge \Phi')$$

For any boolean formula  $\Psi$  and any literal  $x$ , let  $\Psi|x$  (pronounced “sigh given eks”) denote the simpler boolean formula obtained by assuming  $x$  is true. It’s not hard to prove by induction (hint, hint) that  $x \wedge \Psi = x \wedge \Psi|x$ , which implies that

$$\Phi = (x \wedge \Phi'|x) \vee (y \wedge \Phi'|y) \vee (z \wedge \Phi'|z).$$

Thus, in any satisfying assignment for  $\Phi$ , either  $x$  is true and  $\Phi'|x$  is satisfiable, or  $y$  is true and  $\Phi'|y$  is satisfiable, or  $z$  is true and  $\Phi'|z$  is satisfiable. Each of the smaller formulas has at most  $n - 1$  variables. If we recursively check all three possibilities, we get the running time recurrence

$$T(n) \leq 3T(n - 1) + \text{poly}(n),$$

whose solution is  $O(3^n \text{poly}(n))$ . So we’ve actually done *worse!*

But these three recursive cases are not mutually exclusive! If  $\Phi'|x$  is *not* satisfiable, then  $x$  *must* be false in any satisfying assignment for  $\Phi$ . So instead of recursively checking  $\Phi'|y$  in the second step, we can check the even simpler formula  $\Phi'|\bar{x}y$ . Similarly, if  $\Phi'|\bar{x}y$  is not satisfiable, then we know that  $y$  must be false in any satisfying assignment, so we can recursively check  $\Phi'|\bar{x}\bar{y}z$  in the third step.

```
3SAT( $\Phi$ ):
  if  $\Phi = \emptyset$ 
    return TRUE
   $(x \vee y \vee z) \wedge \Phi' \leftarrow \Phi$ 
  if 3SAT( $\Phi|x$ )
    return TRUE
  if 3SAT( $\Phi|\bar{x}y$ )
    return TRUE
  return 3SAT( $\Phi|\bar{x}\bar{y}z$ )
```

The running time off this algorithm obeys the recurrence

$$T(n) = T(n - 1) + T(n - 2) + T(n - 3) + \text{poly}(n),$$

where  $\text{poly}(n)$  denotes the polynomial time required to simplify boolean formulas, handle control flow, move stuff into and out of the recursion stack, and so on. The annihilator method gives us the solution

$$T(n) = O(\lambda^n \text{poly}(n)) = O(1.83928675522^n)$$

where  $\lambda \approx 1.83928675521\dots$  is the largest root of the characteristic polynomial  $r^3 - r^2 - r - 1$ . (Notice that we cleverly eliminated the polynomial noise by increasing the base of the exponent ever so slightly.)

We can improve this algorithm further by eliminating *pure* literals from the formula before recursing. A literal  $x$  is *pure* in if it appears in the formula  $\Phi$  but its negation  $\bar{x}$  does not. It's not hard to prove (hint, hint) that if  $\Phi$  has a satisfying assignment, then it has a satisfying assignment where every pure literal is true. If  $\Phi = (x \vee y \vee z) \wedge \Phi'$  has no pure literals, then some in  $\Phi$  contains the literal  $\bar{x}$ , so we can write

$$\Phi = (x \vee y \vee z) \wedge (\bar{x} \vee u \vee v) \wedge \Phi'$$

for some literals  $u$  and  $v$  (each of which might be  $y$ ,  $\bar{y}$ ,  $z$ , or  $\bar{z}$ ). It follows that the first recursive formula  $\Phi|x$  has contains the clause  $(u \vee v)$ . We can recursively eliminate the variables  $u$  and  $v$  just as we tested the variables  $y$  and  $x$  in the second and third cases of our previous algorithm:

$$\Phi|x = (u \vee v) \wedge \Phi'|x = (u \wedge \Phi'|xu) \vee (v \wedge \Phi'|x\bar{u}v).$$

Here is our new faster algorithm:

```
3SAT( $\Phi$ ):
  if  $\Phi = \emptyset$ 
    return TRUE
  if  $\Phi$  has a pure literal  $x$ 
    return 3SAT( $\Phi|x$ )
   $(x \vee y \vee z) \wedge (\bar{x} \vee u \vee v) \wedge \Phi' \leftarrow \Phi$ 
  if 3SAT( $\Phi|xu$ )
    return TRUE
  if 3SAT( $\Phi|x\bar{u}v$ )
    return TRUE
  if 3SAT( $\Phi|\bar{x}y$ )
    return TRUE
  return 3SAT( $\Phi|\bar{x}\bar{y}z$ )
```

The running time  $T(n)$  of this new algorithm satisfies the recurrence

$$T(n) = 2T(n-2) + 2T(n-3) + \text{poly}(n),$$

and the annihilator method implies that

$$T(n) = O(\mu^n \text{poly}(n)) = \boxed{O(1.76929235425^n)}$$

where  $\mu \approx 1.76929235424\dots$  is the largest root of the characteristic polynomial  $r^3 - 2r - 2$ .

Naturally, this approach can be extended much further; since 1998, at least fifteen different 3SAT algorithms have been published, each improving the running time by a small amount. As of 2010, the fastest deterministic algorithm for 3SAT runs in  $O(1.33334^n)$  time<sup>2</sup>, and the fastest *randomized* algorithm runs in  $O(1.32113^n)$  expected time<sup>3</sup>, but there is good reason to believe that these are *not* the best possible.

<sup>2</sup>Robin A. Moser and Dominik Scheder. A full derandomization of Schöning's  $k$ -SAT algorithm. ArXiv:1008.4067, 2010.

<sup>3</sup>Kazuo Iwama, Kazuhisa Seto, Tadashi Takai, and Suguru Tamaki. Improved randomized algorithms for 3-SAT. To appear in Proc. STACS, 2010.

## 4.2 Maximum Independent Set

Now suppose we are given an undirected graph  $G$  and are asked to find the size of the *largest independent set*, that is, the largest subset of the vertices of  $G$  with no edges between them. Once again, we have an obvious recursive algorithm: Try every subset of nodes, and return the largest subset with no edges. Expressed recursively, the algorithm might look like this.

```
MAXIMUMINDSETSIZE( $G$ ):
    if  $G = \emptyset$ 
        return 0
    else
         $v \leftarrow$  any node in  $G$ 
         $withv \leftarrow 1 + \text{MAXIMUMINDSETSIZE}(G \setminus N(v))$ 
         $withoutv \leftarrow \text{MAXIMUMINDSETSIZE}(G \setminus \{v\})$ 
        return  $\max\{withv, withoutv\}.$ 
```

Here,  $N(v)$  denotes the *neighborhood* of  $v$ : The set containing  $v$  and all of its neighbors. Our algorithm is exploiting the fact that if an independent set contains  $v$ , then by definition it contains none of  $v$ 's neighbors. In the worst case,  $v$  has no neighbors, so  $G \setminus \{v\} = G \setminus N(v)$ . Thus, the running time of this algorithm satisfies the recurrence  $T(n) = 2T(n - 1) + \text{poly}(n) = O(2^n \text{poly}(n))$ . Surprise, surprise.

This algorithm is mirroring a crude recursive upper bound for the number of *maximal* independent sets in a graph; an independent set is maximal if every vertex in  $G$  is either already in the set or a neighbor of a vertex in the set. If the graph is non-empty, then every maximal independent set either includes or excludes each vertex. Thus, the number of maximal independent sets satisfies the recurrence  $M(n) \leq 2M(n - 1)$ , with base case  $M(1) = 1$ . The annihilator method gives us  $M(n) \leq 2^n - 1$ . The only subset that we aren't counting with this upper bound is the empty set!

We can speed up our algorithm by making several careful modifications to avoid the worst case of the running-time recurrence.

- If  $v$  has no neighbors, then  $N(v) = \{v\}$ , and both recursive calls consider a graph with  $n - 1$  nodes. But in this case,  $v$  is in *every* maximal independent set, so one of the recursive calls is redundant. On the other hand, if  $v$  has at least one neighbor, then  $G \setminus N(v)$  has at most  $n - 2$  nodes. So now we have the following recurrence.

$$T(n) \leq O(\text{poly}(n)) + \max \left\{ \begin{array}{l} T(n - 1) \\ T(n - 1) + T(n - 2) \end{array} \right\} = O(1.61803398875^n)$$

The upper bound is derived by solving each case separately using the annihilator method and taking the larger of the two solutions. The first case gives us  $T(n) = O(\text{poly}(n))$ ; the second case yields our old friends the Fibonacci numbers.

- We can improve this bound even more by examining the new worst case:  $v$  has exactly one neighbor  $w$ . In this case, either  $v$  or  $w$  appears in every maximal independent set. However, given any independent set that includes  $w$ , removing  $w$  and adding  $v$  creates another independent set of the same size. It follows that *some maximum independent set includes  $v$* , so we don't need to search the graph  $G \setminus \{v\}$ , and the  $G \setminus N(v)$  has at most  $n - 2$  nodes. On the other hand, if the degree of  $v$  is at least 2, then  $G \setminus N(v)$  has at most  $n - 3$  nodes.

$$T(n) \leq O(\text{poly}(n)) + \max \left\{ \begin{array}{l} T(n - 1) \\ T(n - 2) \\ T(n - 1) + T(n - 3) \end{array} \right\} = O(1.46557123188^n)$$

The base of the exponent is the largest root of the characteristic polynomial  $r^3 - r^2 - 1$ .

- Now the worst-case is a graph where every node has degree at least 2; we split this worst case into two subcases. If  $G$  has a node  $v$  with degree 3 or more, then  $G \setminus N(v)$  has at most  $n - 4$  nodes. Otherwise (since we have already considered nodes of degree 0 and 1), every node in  $G$  has degree 2. Let  $u, v, w$  be a path of three nodes in  $G$  (possibly with  $u$  adjacent to  $w$ ). In any maximal independent set, either  $v$  is present and  $u, w$  are absent, or  $u$  is present and its two neighbors are absent, or  $w$  is present and its two neighbors are absent. In all three cases, we recursively count maximal independent sets in a graph with  $n - 3$  nodes.

$$T(n) \leq O(\text{poly}(n)) + \max \left\{ \begin{array}{l} T(n-1) \\ T(n-2) \\ T(n-1) + T(n-4) \\ 3T(n-3) \end{array} \right\} = O(3^{n/3} \text{poly}(n)) = O(1.44224957031^n)$$

The base of the exponent is  $\sqrt[3]{3}$ , the largest root of the characteristic polynomial  $r^3 - 3$ . The third case would give us a bound of  $O(1.3802775691^n)$ , where the base is the largest root of the characteristic polynomial  $r^4 - r^3 - 1$ .

- Now the worst case for our algorithm is a graph with an extraordinarily special structure: *Every node has degree 2*. In other words, every component of  $G$  is a cycle. But it is easy to prove that the largest independent set in a cycle of length  $k$  has size  $\lfloor k/2 \rfloor$ . So we can handle this case directly in polynomial time, without no recursion at all!

$$T(n) \leq O(\text{poly}(n)) + \max \left\{ \begin{array}{l} T(n-1) \\ T(n-2) \\ T(n-1) + T(n-4) \end{array} \right\} = O(1.3802775691^n)$$

Again, the base of the exponential running time is the largest root of the characteristic polynomial  $r^4 - r^3 - 1$ .

```
MAXIMUMINDSETSIZE( $G$ ):
  if  $G = \emptyset$ 
    return 0
  else if  $G$  has a node  $v$  with degree 0 or 1
    return  $1 + \text{MAXIMUMINDSETSIZE}(G \setminus N(v))$      $\langle\leq n-1\rangle$ 
  else if  $G$  has a node  $v$  with degree greater than 2
     $withv \leftarrow 1 + \text{MAXIMUMINDSETSIZE}(G \setminus N(v))$      $\langle\leq n-4\rangle$ 
     $withoutv \leftarrow \text{MAXIMUMINDSETSIZE}(G \setminus \{v\})$      $\langle\leq n-1\rangle$ 
    return  $\max\{withv, withoutv\}$ 
  else  $\langle\text{every node in } G \text{ has degree 2}\rangle$ 
     $total \leftarrow 0$ 
    for each component of  $G$ 
       $k \leftarrow$  number of vertices in the component
       $total \leftarrow total + \lfloor k/2 \rfloor$ 
    return  $total$ 
```

As with 3SAT, further improvements are possible but increasingly complex. As of 2010, the fastest published algorithm for computing maximum independent sets runs in  $O(1.2210^n)$  time<sup>4</sup>. However,

<sup>4</sup>Fedor V. Fomin, Fabrizio Grandoni, and Dieter Kratsch. Measure and conquer: A simple  $O(2^{0.288n})$  independent set algorithm. Proc. SODA, 18–25, 2006.

in an unpublished technical report, Robson describes a *computer-generated* algorithm that runs in  $O(2^{n/4} \text{poly}(n)) = O(1.1889^n)$  time; just the description of this algorithm requires more than 15 pages.<sup>5</sup>

## Exercises

1. (a) Prove that any  $n$ -vertex graph has at most  $3^{n/3}$  maximal independent sets. [*Hint: Modify the MAXIMUMINDSETSIZE algorithm so that it lists all maximal independent sets.*]
- (b) Describe an  $n$ -vertex graph with exactly  $3^{n/3}$  maximal independent sets, for every integer  $n$  that is a multiple of 3.
- \*2. Describe an algorithm to solve 3SAT in time  $O(\phi^n \text{poly}(n))$ , where  $\phi = (1 + \sqrt{5})/2 \approx 1.618034$ . [*Hint: Prove that in each recursive call, either you have just eliminated a pure literal, or the formula has a clause with at most two literals. What recurrence leads to this running time?*]

---

<sup>5</sup>Mike Robson. Finding a maximum independent set in time  $O(2^{n/4})$ . Technical report 1251-01, LaBRI, 2001. (<http://www.labri.fr/perso/robson/mis/techrep.ps>).

*Those who cannot remember the past are doomed to repeat it.*

— George Santayana, *The Life of Reason, Book I: Introduction and Reason in Common Sense* (1905)

*The 1950s were not good years for mathematical research. We had a very interesting gentleman in Washington named Wilson. He was secretary of Defense, and he actually had a pathological fear and hatred of the word ‘research’. I’m not using the term lightly; I’m using it precisely. His face would suffuse, he would turn red, and he would get violent if people used the term ‘research’ in his presence. You can imagine how he felt, then, about the term ‘mathematical’. The RAND Corporation was employed by the Air Force, and the Air Force had Wilson as its boss, essentially. Hence, I felt I had to do something to shield Wilson and the Air Force from the fact that I was really doing mathematics inside the RAND Corporation. What title, what name, could I choose?*

— Richard Bellman, on the origin of his term ‘dynamic programming’ (1984)

*If we all listened to the professor, we may be all looking for professor jobs.*

— Pittsburgh Steelers’ head coach Bill Cowher, responding to David Romer’s dynamic-programming analysis of football strategy (2003)

## 5 Dynamic Programming

### 5.1 Fibonacci Numbers

#### 5.1.1 Recursive Definitions Are Recursive Algorithms

The Fibonacci numbers  $F_n$ , named after Leonardo Fibonacci Pisano<sup>1</sup>, the mathematician who popularized ‘algorism’ in Europe in the 13th century, are defined as follows:  $F_0 = 0$ ,  $F_1 = 1$ , and  $F_n = F_{n-1} + F_{n-2}$  for all  $n \geq 2$ . The recursive definition of Fibonacci numbers immediately gives us a recursive algorithm for computing them:

```
RECFIBO(n):
    if (n < 2)
        return n
    else
        return RECFIBO(n - 1) + RECFIBO(n - 2)
```

How long does this algorithm take? Except for the recursive calls, the entire algorithm requires only a constant number of steps: one comparison and possibly one addition. If  $T(n)$  represents the number of recursive calls to RECFIBO, we have the recurrence

$$T(0) = 1, \quad T(1) = 1, \quad T(n) = T(n - 1) + T(n - 2) + 1.$$

This looks an awful lot like the recurrence for Fibonacci numbers! The annihilator method gives us an asymptotic bound of  $\Theta(\phi^n)$ , where  $\phi = (\sqrt{5} + 1)/2 \approx 1.61803398875$ , the so-called *golden ratio*, is the largest root of the polynomial  $r^2 - r - 1$ . But it’s fairly easy to prove (hint, hint) the exact solution  $T(n) = 2F_{n+1} - 1$ . In other words, computing  $F_n$  using this algorithm takes more than twice as many steps as just counting to  $F_n$ !

Another way to see this is that the RECFIBO is building a big binary tree of additions, with nothing but zeros and ones at the leaves. Since the eventual output is  $F_n$ , our algorithm must call RECRIBO(1) (which returns 1) exactly  $F_n$  times. A quick inductive argument implies that RECFIBO(0) is called exactly  $F_{n-1}$  times. Thus, the recursion tree has  $F_n + F_{n-1} = F_{n+1}$  leaves, and therefore, because it’s a full binary tree, it must have  $2F_{n+1} - 1$  nodes.

<sup>1</sup>literally, “Leonardo, son of Bonacci, of Pisa”

### 5.1.2 Memo(r)ization: Remember Everything

The obvious reason for the recursive algorithm's lack of speed is that it computes the same Fibonacci numbers over and over and over. A single call to `RECURSIVEFIBO(n)` results in one recursive call to `RECURSIVEFIBO(n - 1)`, two recursive calls to `RECURSIVEFIBO(n - 2)`, three recursive calls to `RECURSIVEFIBO(n - 3)`, five recursive calls to `RECURSIVEFIBO(n - 4)`, and in general,  $F_{k-1}$  recursive calls to `RECURSIVEFIBO(n - k)`, for any  $0 \leq k < n$ . For each call, we're recomputing some Fibonacci number from scratch.

We can speed up our recursive algorithm considerably just by writing down the results of our recursive calls and looking them up again if we need them later. This process was dubbed *memoization* by Richard Michie in the late 1960s.<sup>2</sup>

```
MEMFIBO(n):
  if ( $n < 2$ )
    return  $n$ 
  else
    if  $F[n]$  is undefined
       $F[n] \leftarrow \text{MEMFIBO}(n - 1) + \text{MEMFIBO}(n - 2)$ 
    return  $F[n]$ 
```

Memoization clearly decreases the running time of the algorithm, but by how much? If we actually trace through the recursive calls made by `MEMFIBO`, we find that the array  $F[ ]$  is filled from the bottom up: first  $F[2]$ , then  $F[3]$ , and so on, up to  $F[n]$ . This pattern can be verified by induction: Each entry  $F[i]$  is filled only after its predecessor  $F[i - 1]$ . If we ignore the time spent in recursive calls, it requires only constant time to evaluate the recurrence for each Fibonacci number  $F_i$ . But by design, the recurrence for  $F_i$  is evaluated only once! We conclude that `MEMFIBO` performs only  $O(n)$  additions, an *exponential* improvement over the naïve recursive algorithm!

### 5.1.3 Dynamic Programming: Fill Deliberately

But once we see how the array  $F[ ]$  is filled, we can replace the recursion with a simple loop that intentionally fills the array in order, instead of relying on the complicated recursion to do it for us ‘accidentally’.

```
ITERFIBO(n):
   $F[0] \leftarrow 0$ 
   $F[1] \leftarrow 1$ 
  for  $i \leftarrow 2$  to  $n$ 
     $F[i] \leftarrow F[i - 1] + F[i - 2]$ 
  return  $F[n]$ 
```

Now the time analysis is immediate: `ITERFIBO` clearly uses  **$O(n)$  additions** and stores  **$O(n)$  integers**.

This gives us our first explicit *dynamic programming* algorithm. The dynamic programming paradigm was developed by Richard Bellman in the mid-1950s, while working at the RAND Corporation. Bellman deliberately chose the name ‘dynamic programming’ to hide the mathematical character of his work from his military bosses, who were actively hostile toward anything resembling mathematical research. Here, the word ‘programming’ does not refer to writing code, but rather to the older sense of *planning* or *scheduling*, typically by filling in a table. For example, sports programs and theater programs are schedules of important events (with ads); television programming involves filling each available time slot with a show (and ads); degree programs are schedules of classes to be taken (with ads). The Air

<sup>2</sup>“My name is Elmer J. Fudd, millionaire. I own a mansion and a yacht.”

Force funded Bellman and others to develop methods for constructing training and logistics schedules, or as they called them, ‘programs’. The word ‘dynamic’ is meant to suggest that the table is filled in over time, rather than all at once (as in ‘linear programming’, which we will see later in the semester).<sup>3</sup>

### 5.1.4 Don’t Remember Everything After All

In many dynamic programming algorithms, it is not necessary to retain *all* intermediate results through the entire computation. For example, we can significantly reduce the space requirements of our algorithm ITERFIBO by maintaining only the two newest elements of the array:

```
ITERFIBO2( $n$ ):
    prev ← 1
    curr ← 0
    for  $i \leftarrow 1$  to  $n$ 
        next ← curr + prev
        prev ← curr
        curr ← next
    return curr
```

(This algorithm uses the non-standard but perfectly consistent base case  $F_{-1} = 1$  so that ITERFIBO2(0) returns the correct value 0.)

### 5.1.5 Faster! Faster!

Even this algorithm can be improved further, using the following wonderful fact:

$$\begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} y \\ x+y \end{bmatrix}$$

In other words, multiplying a two-dimensional vector by the matrix  $\begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}$  does exactly the same thing as one iteration of the inner loop of ITERFIBO2. This might lead us to believe that multiplying by the matrix  $n$  times is the same as iterating the loop  $n$  times:

$$\begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}^n \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} F_{n-1} \\ F_n \end{bmatrix}.$$

A quick inductive argument proves this fact. So if we want the  $n$ th Fibonacci number, we just have to compute the  $n$ th power of the matrix  $\begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}$ . If we use repeated squaring, computing the  $n$ th power of something requires only  $O(\log n)$  multiplications. In this case, that means  $O(\log n)$   $2 \times 2$  matrix multiplications, each of which reduces to a constant number of integer multiplications and additions. Thus, we can compute  $F_n$  in only  **$O(\log n)$  integer arithmetic operations**.

This is an exponential speedup over the standard iterative algorithm, which was already an exponential speedup over our original recursive algorithm. Right?

### 5.1.6 Whoa! Not so fast!

Well, not exactly. Fibonacci numbers grow exponentially fast. The  $n$ th Fibonacci number is approximately  $n \log_{10} \phi \approx n/5$  decimal digits long, or  $n \log_2 \phi \approx 2n/3$  bits. So we can’t possibly compute  $F_n$  in logarithmic time — we need  $\Omega(n)$  time just to write down the answer!

---

<sup>3</sup>“I thought dynamic programming was a good name. It was something not even a Congressman could object to. So I used it as an umbrella for my activities.”

The way out of this apparent paradox is to observe that *we can't perform arbitrary-precision arithmetic in constant time*. Multiplying two  $n$ -digit numbers using fast Fourier transforms (described in a different lecture note) requires  $O(n \log n \log \log n)$  time. Thus, the matrix-based algorithm's actual running time obeys the recurrence  $T(n) = T(\lfloor n/2 \rfloor) + O(n \log n \log \log n)$ , which solves to  $T(n) = O(n \log n \log \log n)$  using recursion trees.

Is this slower than our “linear-time” iterative algorithm? No! Addition isn’t free, either. Adding two  $n$ -digit numbers takes  $O(n)$  time, so the running time of the iterative algorithm is  $O(n^2)$ . (Do you see why?) Our matrix algorithm really is faster than our iterative algorithm, but not exponentially faster.

In the original recursive algorithm, the extra cost of arbitrary-precision arithmetic is overwhelmed by the huge number of recursive calls. The correct recurrence is  $T(n) = T(n - 1) + T(n - 2) + O(n)$ , for which the annihilator method still implies the solution  $T(n) = O(\phi^n)$ .

## 5.2 Longest Increasing Subsequence

In a previous lecture, we developed a recursive algorithm to find the length of the longest increasing subsequence of a given sequence of numbers. Given an array  $A[1..n]$ , the length of the longest increasing subsequences is computed by the function call  $\text{LISBIGGER}(-\infty, A[1..n])$ , where  $\text{LISBIGGER}$  is the following recursive algorithm:

```
LISBIGGER(prev, A[1..n]):
    if  $n = 0$ 
        return 0
    else
         $max \leftarrow \text{LISBIGGER}(prev, A[2..n])$ 
        if  $A[1] > prev$ 
             $L \leftarrow 1 + \text{LISBIGGER}(A[1], A[2..n])$ 
            if  $L > max$ 
                 $max \leftarrow L$ 
        return  $max$ 
```

We can simplify our notation slightly with two simple observations. First, the input variable  $prev$  is always either  $-\infty$  or an element of the input array. Second, the second argument of  $\text{LISBIGGER}$  is always a *suffix* of the original input array. If we add a new sentinel value  $A[0] = -\infty$  to the input array, we can identify any recursive subproblem with two array indices.

Thus, we can rewrite the recursive algorithm as follows. Add the sentinel value  $A[0] = -\infty$ . Let  $LIS(i, j)$  denote the length of the longest increasing subsequence of  $A[j..n]$  with all elements larger than  $A[i]$ . Our goal is to compute  $LIS(0, 1)$ . For all  $i < j$ , we have

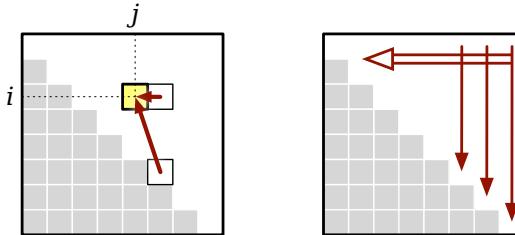
$$LIS(i, j) = \begin{cases} 0 & \text{if } j > n \\ LIS(i, j + 1) & \text{if } A[i] \geq A[j] \\ \max\{LIS(i, j + 1), 1 + LIS(j, j + 1)\} & \text{otherwise} \end{cases}$$

Because each recursive subproblem can be identified by two indices  $i$  and  $j$ , we can store the intermediate values in a two-dimensional array  $LIS[0..n, 1..n]$ . Since there are  $O(n^2)$  entries in the table, our memoized algorithm uses  $O(n^2)$  space.<sup>4</sup> Each entry in the table can be computed in  $O(1)$  time once we know its predecessors, so our memoized algorithm runs in  $O(n^2)$  time.

---

<sup>4</sup>In fact, we only need half of this array, because we always have  $i < j$ . But even if we cared about constant factors in this class (we don’t), this would be the wrong time to worry about them. The first order of business is to find an algorithm that actually works; once we have that, then we can think about optimizing it.

It's not immediately clear what order the recursive algorithm fills the rest of the table; all we can tell from the recurrence is that each entry  $LIS[i, j]$  is filled in *after* the entries  $LIS[i, j+1]$  and  $LIS[j, j+1]$  in the next columns. But just this partial information is enough to give us an explicit evaluation order. If we fill in our table one column at a time, from right to left, then whenever we reach an entry in the table, the entries it depends on are already available.



Dependencies in the memoization table for longest increasing subsequence, and a legal evaluation order

Finally, putting everything together, we obtain the following dynamic programming algorithm:

```

LIS(A[1 .. n]):
   $A[0] \leftarrow -\infty$            «Add a sentinel»
  for  $i \leftarrow 0$  to  $n$           «Base cases»
     $LIS[i, n+1] \leftarrow 0$ 
  for  $j \leftarrow n$  downto 1
    for  $i \leftarrow 0$  to  $j-1$ 
      if  $A[i] \geq A[j]$ 
         $LIS[i, j] \leftarrow LIS[i, j+1]$ 
      else
         $LIS[i, j] \leftarrow \max\{LIS[i, j+1], 1 + LIS[j, j+1]\}$ 
  return  $LIS[0, 1]$ 

```

As expected, the algorithm clearly uses  $O(n^2)$  time and space. However, we can reduce the space to  $O(n)$  by only maintaining the two most recent columns of the table,  $LIS[\cdot, j]$  and  $LIS[\cdot, j+1]$ .<sup>5</sup>

This is not the only recursive strategy we could use for computing longest increasing subsequences efficiently. Here is another recurrence that gives us the  $O(n)$  space bound for free. Let  $LIS'(i)$  denote the length of the longest increasing subsequence of  $A[i..n]$  that starts with  $A[i]$ . Our goal is to compute  $LIS'(0) - 1$ ; we subtract 1 to ignore the sentinel value  $-\infty$ . To define  $LIS'(i)$  recursively, we only need to specify the second element in subsequence; the Recursion Fairy will do the rest.

$$LIS'(i) = 1 + \max \{LIS'(j) \mid j > i \text{ and } A[j] > A[i]\}$$

Here, I'm assuming that  $\max \emptyset = 0$ , so that the base case is  $L'(n) = 1$  falls out of the recurrence automatically. Memoizing this recurrence requires only  $O(n)$  space, and the resulting algorithm runs in  $O(n^2)$  time. To transform this memoized recurrence into a dynamic programming algorithm, we only need to guarantee that  $LIS'(j)$  is computed before  $LIS'(i)$  whenever  $i < j$ .

<sup>5</sup>See, I told you not to worry about constant factors yet!

```

LIS2( $A[1..n]$ ):
 $A[0] = -\infty$             $\langle\langle$  Add a sentinel  $\rangle\rangle$ 
for  $i \leftarrow n$  downto 0
   $LIS'[i] \leftarrow 1$ 
  for  $j \leftarrow i + 1$  to  $n$ 
    if  $A[j] > A[i]$  and  $1 + LIS'[j] > LIS'[i]$ 
       $LIS'[i] \leftarrow 1 + LIS'[j]$ 
return  $LIS'[0] - 1$          $\langle\langle$  Don't count the sentinel  $\rangle\rangle$ 

```

### 5.3 The Pattern: Smart Recursion

In a nutshell, dynamic programming is *recursion without repetition*. Dynamic programming algorithms store the solutions of intermediate subproblems, often *but not always* in some kind of array or table. Many algorithms students make the mistake of focusing on the table (because tables are easy and familiar) instead of the *much* more important (and difficult) task of finding a correct recurrence. As long as we memoize the correct recurrence, an explicit table isn't really necessary, but if the recursion is incorrect, nothing works.

**Dynamic programming is *not* about filling in tables.  
It's about smart recursion!**

Dynamic programming algorithms are almost always developed in two distinct stages.

1. **Formulate the problem recursively.** Write down a recursive formula or algorithm for the whole problem in terms of the answers to smaller subproblems. This is the hard part.
2. **Build solutions to your recurrence from the bottom up.** Write an algorithm that starts with the base cases of your recurrence and works its way up to the final solution, by considering intermediate subproblems in the correct order. This stage can be broken down into several smaller, relatively mechanical steps:
  - (a) **Identify the subproblems.** What are all the different ways can your recursive algorithm call itself, starting with some initial input? For example, the argument to REC<sub>FIBO</sub> is always an integer between 0 and  $n$ .
  - (b) **Analyze space and running time.** The number of possible distinct subproblems determines the space complexity of your memoized algorithm. To compute the time complexity, add up the running times of all possible subproblems, *ignoring the recursive calls*. For example, if we already know  $F_{i-1}$  and  $F_{i-2}$ , we can compute  $F_i$  in  $O(1)$  time, so computing the first  $n$  Fibonacci numbers takes  $O(n)$  time.
  - (c) **Choose a data structure to memoize intermediate results.** For most problems, each recursive subproblem can be identified by a few integers, so you can use a multidimensional array. For some problems, however, a more complicated data structure is required.
  - (d) **Identify dependencies between subproblems.** Except for the base cases, every recursive subproblem depends on other subproblems—which ones? Draw a picture of your data structure, pick a generic element, and draw arrows from each of the other elements it depends on. Then formalize your picture.
  - (e) **Find a good evaluation order.** Order the subproblems so that each subproblem comes *after* the subproblems it depends on. Typically, this means you should consider the base cases

first, then the subproblems that depends only on base cases, and so on. More formally, the dependencies you identified in the previous step define a partial order over the subproblems; in this step, you need to find a linear extension of that partial order. *Be careful!*

- (f) **Write down the algorithm.** You know what order to consider the subproblems, and you know how to solve each subproblem. So do that! If your data structure is an array, this usually means writing a few nested for-loops around your original recurrence. **You don't need to do this on homework or exams.**

Of course, you have to prove that each of these steps is correct. If your recurrence is wrong, or if you try to build up answers in the wrong order, your algorithm won't work!

#### 5.4 Warning: Greed is Stupid

If we're very very very lucky, we can bypass all the recurrences and tables and so forth, and solve the problem using a *greedy* algorithm. The general greedy strategy is look for the best first step, take it, and then continue. While this approach seems very natural, it almost never works; optimization problems that can be solved correctly by a greedy algorithm are very rare. Nevertheless, for many problems that should be solved by dynamic programming, many students' first intuition is to apply a greedy strategy.

For example, a greedy algorithm for the edit distance problem might look for the longest common substring of the two strings, match up those substrings (since those substitutions don't cost anything), and then recursively look for the edit distances between the left halves and right halves of the strings. If there is no common substring—that is, if the two strings have no characters in common—the edit distance is clearly the length of the larger string. If this sounds like a stupid hack to you, pat yourself on the back. It isn't even *close* to the correct solution.

Everyone should tattoo the following sentence on the back of their hands, right under all the rules about logarithms and big-Oh notation:

**Greedy algorithms never work!  
Use dynamic programming instead!**

What, never?

No, never!

What, *never*?

Well... hardly ever.<sup>6</sup>

A different lecture note describes the effort required to prove that greedy algorithms are correct, in the rare instances when they are. **You will not receive any credit for any greedy algorithm for any problem in this class without a formal proof of correctness.** We'll push through the formal proofs for several greedy algorithms later in the semester.

#### 5.5 Edit Distance

The *edit distance* between two words—sometimes also called the *Levenshtein distance*—is the minimum number of letter insertions, letter deletions, and letter substitutions required to transform one word into another. For example, the edit distance between FOOD and MONEY is at most four:

---

<sup>6</sup>Greedy methods hardly ever work! So give three cheers, and one cheer more, for the hardy Captain of the *Pinafore!* Then give three cheers, and one cheer more, for the Captain of the *Pinafore!*

FOOD → MOD → MOND → MONE → MONEY

A better way to display this editing process is to place the words one above the other, with a gap in the first word for every insertion, and a gap in the second word for every deletion. Columns with two *different* characters correspond to substitutions. Thus, the number of editing steps is just the number of columns that don't contain the same character twice.

F	O	O	D	
M	O	N	E	Y

It's fairly obvious that you can't get from FOOD to MONEY in three steps, so their edit distance is exactly four. Unfortunately, this is not so easy in general. Here's a longer example, showing that the distance between ALGORITHM and ALTRUISTIC is at most six. Is this optimal?

A	L	G	O	R	I	T	H	M		
A	L		T	R	U	I	S	T	I	C

To develop a dynamic programming algorithm to compute the edit distance between two strings, we first need to develop a recursive definition. Our gap representation for edit sequences has a crucial “optimal substructure” property. Suppose we have the gap representation for the shortest edit sequence for two strings. **If we remove the last column, the remaining columns must represent the shortest edit sequence for the remaining substrings.** We can easily prove this by contradiction. If the substrings had a shorter edit sequence, we could just glue the last column back on and get a shorter edit sequence for the original strings. Once we figure out what should go in the last column, the Recursion Fairy will magically give us the rest of the optimal gap representation.

So let's recursively define the edit distance between two strings  $A[1..m]$  and  $B[1..n]$ , which we denote by  $Edit(A[1..m], B[1..n])$ . If neither string is empty, there are three possibilities for the last column in the shortest edit sequence:

- **Insertion:** The last entry in the bottom row is empty. In this case, the edit distance is equal to  $Edit(A[1..m-1], B[1..n]) + 1$ . The  $+1$  is the cost of the final insertion, and the recursive expression gives the minimum cost for the other columns.
- **Deletion:** The last entry in the top row is empty. In this case, the edit distance is equal to  $Edit(A[1..m], B[1..n-1]) + 1$ . The  $+1$  is the cost of the final deletion, and the recursive expression gives the minimum cost for the other columns.
- **Substitution:** Both rows have characters in the last column. If the characters are the same, the substitution is free, so the edit distance is equal to  $Edit(A[1..m-1], B[1..n-1])$ . If the characters are different, then the edit distance is equal to  $Edit(A[1..m-1], B[1..n-1]) + 1$ .

The edit distance between  $A$  and  $B$  is the smallest of these three possibilities:<sup>7</sup>

$$Edit(A[1..m], B[1..n]) = \min \left\{ \begin{array}{l} Edit(A[1..m-1], B[1..n]) + 1 \\ Edit(A[1..m], B[1..n-1]) + 1 \\ Edit(A[1..m-1], B[1..n-1]) + [A[m] \neq B[n]] \end{array} \right\}$$

<sup>7</sup>Once again, I'm using Iverson's bracket notation  $[P]$  to denote the *indicator variable* for the logical proposition  $P$ , which has value 1 if  $P$  is true and 0 if  $P$  is false.

This recurrence has two easy base cases. The only way to convert the empty string into a string of  $n$  characters is by performing  $n$  insertions. Similarly, the only way to convert a string of  $m$  characters into the empty string is with  $m$  deletions. Thus, if  $\epsilon$  denotes the empty string, we have

$$\text{Edit}(A[1..m], \epsilon) = m, \quad \text{Edit}(\epsilon, B[1..n]) = n.$$

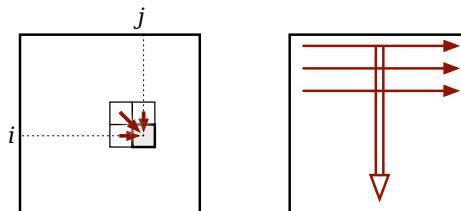
Both of these expressions imply the trivial base case  $\text{Edit}(\epsilon, \epsilon) = 0$ .

Now notice that the arguments to our recursive subproblems are always *prefixes* of the original strings  $A$  and  $B$ . Thus, we can simplify our notation considerably by using the lengths of the prefixes, instead of the prefixes themselves, as the arguments to our recursive function. So let's write  $\text{Edit}(i, j)$  as shorthand for  $\text{Edit}(A[1..i], B[1..j])$ . This function satisfies the following recurrence:

$$\boxed{\text{Edit}(i, j) = \begin{cases} i & \text{if } j = 0 \\ j & \text{if } i = 0 \\ \min \left\{ \begin{array}{l} \text{Edit}(i - 1, j) + 1, \\ \text{Edit}(i, j - 1) + 1, \\ \text{Edit}(i - 1, j - 1) + [A[i] \neq B[j]] \end{array} \right\} & \text{otherwise} \end{cases}}$$

The edit distance between the original strings  $A$  and  $B$  is just  $\text{Edit}(m, n)$ . This recurrence translates directly into a recursive algorithm; the precise running time is not obvious, but it's clearly exponential in  $m$  and  $n$ . **Fortunately, we don't care about the precise running time of the recursive algorithm.** The recursive running time wouldn't tell us anything about our eventual dynamic programming algorithm, so we're just not going to bother computing it.<sup>8</sup>

Because each recursive subproblem can be identified by two indices  $i$  and  $j$ , we can memoize intermediate values in a two-dimensional array  $\text{Edit}[0..m, 0..n]$ . Note that the index ranges start at zero to accommodate the base cases. Since there are  $\Theta(mn)$  entries in the table, our memoized algorithm uses  **$\Theta(mn)$  space**. Since each entry in the table can be computed in  $\Theta(1)$  time once we know its predecessors, our memoized algorithm runs in  **$\Theta(mn)$  time**.



Dependencies in the memoization table for edit distance, and a legal evaluation order

<sup>8</sup>In case you're curious, the running time of the unmemoized recursive algorithm obeys the following recurrence:

$$T(m, n) = \begin{cases} O(1) & \text{if } n = 0 \text{ or } m = 0, \\ T(m, n - 1) + T(m - 1, n) + T(n - 1, m - 1) + O(1) & \text{otherwise.} \end{cases}$$

I don't know of a general closed-form solution for this mess, but we can derive an upper bound by defining a new function

$$T'(N) = \max_{n+m=N} T(n, m) = \begin{cases} O(1) & \text{if } N = 0, \\ 2T(N - 1) + T(N - 2) + O(1) & \text{otherwise.} \end{cases}$$

The annihilator method implies that  $T'(N) = O((1 + \sqrt{2})^N)$ . Thus, the running time of our recursive edit-distance algorithm is at most  $T'(n + m) = O((1 + \sqrt{2})^{n+m})$ .

Each entry  $Edit[i, j]$  depends only on its three neighboring entries  $Edit[i - 1, j] + 1$ ,  $Edit[i, j - 1] + 1$ , and  $Edit[i - 1, j - 1]$ . If we fill in our table in the standard row-major order—row by row from top down, each row from left to right—then whenever we reach an entry in the table, the entries it depends on are already available. Putting everything together, we obtain the following dynamic programming algorithm:

```
EDITDISTANCE( $A[1..m], B[1..n]$ ):
    for  $j \leftarrow 1$  to  $n$ 
         $Edit[0, j] \leftarrow j$ 
    for  $i \leftarrow 1$  to  $m$ 
         $Edit[i, 0] \leftarrow i$ 
    for  $j \leftarrow 1$  to  $n$ 
        if  $A[i] = B[j]$ 
             $Edit[i, j] \leftarrow \min\{Edit[i - 1, j] + 1, Edit[i, j - 1] + 1, Edit[i - 1, j - 1]\}$ 
        else
             $Edit[i, j] \leftarrow \min\{Edit[i - 1, j] + 1, Edit[i, j - 1] + 1, Edit[i - 1, j - 1] + 1\}$ 
    return  $Edit[m, n]$ 
```

Here's the resulting table for ALGORITHM → ALTRUISTIC. Bold numbers indicate places where characters in the two strings are equal. The arrows represent the predecessor(s) that actually define each entry. Each direction of arrow corresponds to a different edit operation: horizontal=deletion, vertical=insertion, and diagonal=substitution. Bold diagonal arrows indicate “free” substitutions of a letter for itself. Any path of arrows from the top left corner to the bottom right corner of this table represents an optimal edit sequence between the two strings. (There can be many such paths.) Moreover, since we can compute these arrows in a post-processing phase from the values stored in the table, we can reconstruct the actual optimal editing sequence in  $O(n + m)$  additional time.

	A	L	G	O	R	I	T	H	M
	0 → 1 → 2 → 3 → 4 → 5 → 6 → 7 → 8 → 9								
A	1	<b>0</b> → 1 → 2 → 3 → 4 → 5 → 6 → 7 → 8							
L	2	1	<b>0</b> → 1 → 2 → 3 → 4 → 5 → 6 → 7						
T	3	2	1	1 → 2 → 3 → 4 → <b>5</b> → 6 → 7					
R	4	3	2	2	2 → 3 → 4 → <b>5</b> → 6				
U	5	4	3	3	3	3 → 4 → 5 → 6			
I	6	5	4	4	4	<b>3</b> → 4 → 5 → 6			
S	7	6	5	5	5	4	4	5	6
T	8	7	6	6	6	5	<b>4</b> → 5 → 6		
I	9	8	7	7	7	7	<b>6</b>	5	5 → 6
C	10	9	8	8	8	8	7	6	6 → 6

The edit distance between ALGORITHM and ALTRUISTIC is indeed six. There are three paths through this table from the top left to the bottom right, so there are three optimal edit sequences:

A	L	G	O	R	I	T	H	M	
A	L	T	R	U	I	S	T	I	C

A	L	G	O	R	I	T	H	M		
A	L		T	R	U	I	S	T	I	C

A	L	G	O	R	I	T	H	M	
A	L	T	R	U	I	S	T	I	C

## 5.6 Optimal Binary Search Trees

In an earlier lecture, we developed a recursive algorithm for the optimal binary search tree problem. We are given a sorted array  $A[1..n]$  of search keys and an array  $f[1..n]$  of frequency counts, where  $f[i]$  is the number of searches to  $A[i]$ . Our task is to construct a binary search tree for that set such that the total cost of all the searches is as small as possible. We developed the following recurrence for this problem:

$$OptCost(f[1..n]) = \min_{1 \leq r \leq n} \left\{ OptCost(f[1..r-1]) + \sum_{i=1}^n f[i] + OptCost(f[r+1..n]) \right\}$$

To put this recurrence in more standard form, fix the frequency array  $f$ , and let  $OptCost(i, j)$  denote the total search time in the optimal search tree for the subarray  $A[i..j]$ . To simplify notation a bit, let  $F(i, j)$  denote the total frequency count for all the keys in the interval  $A[i..j]$ :

$$F(i, j) := \sum_{k=i}^j f[k]$$

We can now write

$$OptCost(i, j) = \begin{cases} 0 & \text{if } j < i \\ F(i, j) + \min_{i \leq r \leq j} (OptCost(i, r-1) + OptCost(r+1, j)) & \text{otherwise} \end{cases}$$

The base case might look a little weird, but all it means is that the total cost for searching an empty set of keys is zero.

The algorithm will be somewhat simpler and more efficient if we precompute all possible values of  $F(i, j)$  and store them in an array. Computing each value  $F(i, j)$  using a separate for-loop would  $O(n^3)$  time. A better approach is to turn the recurrence

$$F(i, j) = \begin{cases} f[i] & \text{if } i = j \\ F(i, j-1) + f[j] & \text{otherwise} \end{cases}$$

into the following  $O(n^2)$ -time dynamic programming algorithm:

```

INITF( $f[1..n]$ ):
  for  $i \leftarrow 1$  to  $n$ 
     $F[i, i-1] \leftarrow 0$ 
    for  $j \leftarrow i$  to  $n$ 
       $F[i, j] \leftarrow F[i, j-1] + f[i]$ 

```

This will be used as an initialization subroutine in our final algorithm.

So now let's compute the optimal search tree cost  $OptCost(1, n)$  from the bottom up. We can store all intermediate results in a table  $OptCost[1..n, 0..n]$ . Only the entries  $OptCost[i, j]$  with  $j \geq i - 1$  will actually be used. The base case of the recurrence tells us that any entry of the form  $OptCost[i, i - 1]$  can immediately be set to 0. For any other entry  $OptCost[i, j]$ , we can use the following algorithm fragment, which comes directly from the recurrence:

```
COMPUTEOPTCOST( $i, j$ ):
 $OptCost[i, j] \leftarrow \infty$ 
for  $r \leftarrow i$  to  $j$ 
     $tmp \leftarrow OptCost[i, r - 1] + OptCost[r + 1, j]$ 
    if  $OptCost[i, j] > tmp$ 
         $OptCost[i, j] \leftarrow tmp$ 
 $OptCost[i, j] \leftarrow OptCost[i, j] + F[i, j]$ 
```

The only question left is what order to fill in the table.

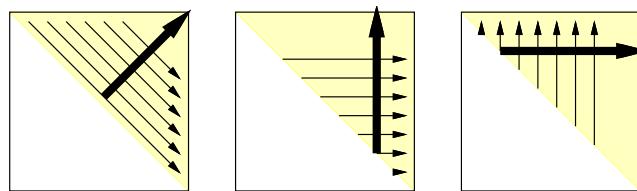
Each entry  $OptCost[i, j]$  depends on all entries  $OptCost[i, r - 1]$  and  $OptCost[r + 1, j]$  with  $i \leq k \leq j$ . In other words, every entry in the table depends on all the entries directly to the left or directly below. In order to fill the table efficiently, we must choose an order that computes all those entries before  $OptCost[i, j]$ . There are at least three different orders that satisfy this constraint. The one that occurs to most people first is to scan through the table one diagonal at a time, starting with the trivial base cases  $OptCost[i, i - 1]$ . The complete algorithm looks like this:

```
OPTIMALSEARCHTREE( $f[1..n]$ ):
INITF( $f[1..n]$ )
for  $i \leftarrow 1$  to  $n$ 
     $OptCost[i, i - 1] \leftarrow 0$ 
for  $d \leftarrow 0$  to  $n - 1$ 
    for  $i \leftarrow 1$  to  $n - d$ 
        COMPUTEOPTCOST( $i, i + d$ )
return  $OptCost[1, n]$ 
```

We could also traverse the array row by row from the bottom up, traversing each row from left to right, or column by column from left to right, traversing each columns from the bottom up. These two orders give us the following algorithms:

```
OPTIMALSEARCHTREE2( $f[1..n]$ ):
INITF( $f[1..n]$ )
for  $i \leftarrow n$  downto 1
     $OptCost[i, i - 1] \leftarrow 0$ 
    for  $j \leftarrow i$  to  $n$ 
        COMPUTEOPTCOST( $i, j$ )
return  $OptCost[1, n]$ 
```

```
OPTIMALSEARCHTREE3( $f[1..n]$ ):
INITF( $f[1..n]$ )
for  $j \leftarrow 0$  to  $n$ 
     $OptCost[j + 1, j] \leftarrow 0$ 
    for  $i \leftarrow j$  downto 1
        COMPUTEOPTCOST( $i, j$ )
return  $OptCost[1, n]$ 
```



Three different evaluation orders for the table  $OptCost[i, j]$ .

No matter which of these orders we actually use, the resulting algorithm runs in  $\Theta(n^3)$  time and uses  $\Theta(n^2)$  space.

We could have predicted these space and time bounds directly from the original recurrence.

$$OptCost(i, j) = \begin{cases} 0 & \text{if } j = i - i \\ F(i, j) + \min_{i \leq r \leq j} (OptCost(i, r - 1) + OptCost(r + 1, j)) & \text{otherwise} \end{cases}$$

First, the function has two arguments, each of which can take on any value between 1 and  $n$ , so we probably need a table of size  $O(n^2)$ . Next, there are three variables in the recurrence ( $i$ ,  $j$ , and  $r$ ), each of which can take any value between 1 and  $n$ , so it should take us  $O(n^3)$  time to fill the table.

## 5.7 Dynamic Programming on Trees

So far, all of our dynamic programming examples use a multidimensional array to store the results of recursive subproblems. However, as the next example shows, this is not always the most appropriate data structure to use.

A **independent set** in a graph is a subset of the vertices that have no edges between them. Finding the largest independent set in an arbitrary graph is extremely hard; in fact, this is one of the canonical NP-hard problems described in another lecture note. But from some special cases of graphs, we can find the largest independent set efficiently. In particular, when the input graph is a tree (a connected and acyclic graph) with  $n$  vertices, we can compute the largest independent set in  $O(n)$  time.

In the recursion notes, we saw a recursive algorithm for computing the size of the largest independent set in an arbitrary graph:

```
MAXIMUMINDSETSIZE( $G$ ):
    if  $G = \emptyset$ 
        return 0
     $v \leftarrow$  any node in  $G$ 
     $withv \leftarrow 1 + \text{MAXIMUMINDSETSIZE}(G \setminus N(v))$ 
     $withoutv \leftarrow \text{MAXIMUMINDSETSIZE}(G \setminus \{v\})$ 
    return  $\max\{withv, withoutv\}$ .
```

Here,  $N(v)$  denotes the *neighborhood* of  $v$ : the set containing  $v$  and all of its neighbors. As we observed in the other lecture notes, this algorithm has a worst-case running time of  $O(2^n \text{poly}(n))$ , where  $n$  is the number of vertices in the input graph.

Now suppose we require that the input graph is a tree; we will call this tree  $T$  instead of  $G$  from now on. We need to make a slight change to the algorithm to make it truly recursive. The subgraphs  $T \setminus \{v\}$  and  $T \setminus N(v)$  are forests, which may have more than one component. But the largest independent set in a disconnected graph is just the union of the largest independent sets in its components, so we can separately consider each tree in these forests. Fortunately, this has the added benefit of making the recursive algorithm more efficient, especially if we can choose the node  $v$  such that the trees are all significantly smaller than  $T$ . Here is the modified algorithm:

```

MAXIMUMINDSETSIZE( $T$ ):
  if  $T = \emptyset$ 
    return 0
   $v \leftarrow$  any node in  $T$ 
   $withv \leftarrow 1$ 
  for each tree  $T'$  in  $T \setminus N(v)$ 
     $withv \leftarrow withv + \text{MAXIMUMINDSETSIZE}(T')$ 
   $withoutv \leftarrow 0$ 
  for each tree  $T'$  in  $T \setminus \{v\}$ 
     $withoutv \leftarrow withoutv + \text{MAXIMUMINDSETSIZE}(T')$ 
  return  $\max\{withv, withoutv\}.$ 

```

Now let's try to memoize this algorithm. Each recursive subproblem considers a subtree (that is, a connected subgraph) of the original tree  $T$ . Unfortunately, a single tree  $T$  can have exponentially many subtrees, so we seem to be doomed from the start!

Fortunately, there's a degree of freedom that we have not yet exploited: *We get to choose the vertex  $v$ .* We need a recipe—an algorithm!—for choosing  $v$  in each subproblem that limits the number of different subproblems the algorithm considers. To make this work, we impose some additional structure on the original input tree. Specifically, we declare one of the vertices of  $T$  to be the *root*, and we orient all the edges of  $T$  away from that root. Then we let  $v$  be the root of the input tree; this choice guarantees that each recursive subproblem considers a *rooted* subtree of  $T$ . Each vertex in  $T$  is the root of exactly one subtree, so now the number of distinct subproblems is exactly  $n$ . We can further simplify the algorithm by only passing a single node instead of the entire subtree:

```

MAXIMUMINDSETSIZE( $v$ ):
   $withv \leftarrow 1$ 
  for each grandchild  $x$  of  $v$ 
     $withv \leftarrow withv + \text{MAXIMUMINDSETSIZE}(x)$ 
   $withoutv \leftarrow 0$ 
  for each child  $w$  of  $v$ 
     $withoutv \leftarrow withoutv + \text{MAXIMUMINDSETSIZE}(w)$ 
  return  $\max\{withv, withoutv\}.$ 

```

What data structure should we use to store intermediate results? The most natural choice is the tree itself! Specifically, for each node  $v$ , we store the result of  $\text{MAXIMUMINDSETSIZE}(v)$  in a new field  $v.MIS$ . (We *could* use an array, but then we'd have to add a new field to each node anyway, pointing to the corresponding array entry. Why bother?)

What's the running time of the algorithm? The non-recursive time associated with each node  $v$  is proportional to the number of children and grandchildren of  $v$ ; this number can be very different from one vertex to the next. But we can turn the analysis around: Each vertex contributes a constant amount of time to its parent and its grandparent! Since each vertex has at most one parent and at most one grandparent, the total running time is  $O(n)$ .

What's a good order to consider the subproblems? The subproblem associated with any node  $v$  depends on the subproblems associated with the children and grandchildren of  $v$ . So we can visit the nodes in any order, provided that all children are visited before their parent. In particular, we can use a straightforward post-order traversal.

Here is the resulting dynamic programming algorithm. Yes, it's still recursive. I've swapped the evaluation of the with- $v$  and without- $v$  cases; we need to visit the kids first anyway, so why not consider the subproblem that depends directly on the kids first?

```

MAXIMUMINDSETSIZE( $v$ ):
     $withoutv \leftarrow 0$ 
    for each child  $w$  of  $v$ 
         $withoutv \leftarrow withoutv + \text{MAXIMUMINDSETSIZE}(w)$ 
     $withv \leftarrow 1$ 
    for each grandchild  $x$  of  $v$ 
         $withv \leftarrow withv + x.\text{MIS}$ 
     $v.\text{MIS} \leftarrow \max\{withv, withoutv\}$ 
    return  $v.\text{MIS}$ 

```

Another option is to store *two* values for each rooted subtree: the size of the largest independent set *that includes the root*, and the size of the largest independent set *that excludes the root*. This gives us an even simpler algorithm, with the same  $O(n)$  running time.

```

MAXIMUMINDSETSIZE( $v$ ):
     $v.\text{MISno} \leftarrow 0$ 
     $v.\text{MISyes} \leftarrow 1$ 
    for each child  $w$  of  $v$ 
         $v.\text{MISno} \leftarrow v.\text{MISno} + \text{MAXIMUMINDSETSIZE}(w)$ 
         $v.\text{MISyes} \leftarrow v.\text{MISyes} + w.\text{MISno}$ 
    return  $\max\{v.\text{MISyes}, v.\text{MISno}\}$ 

```

## Exercises

1. Suppose you are given an array  $A[1..n]$  of numbers, which may be positive, negative, or zero.
  - (a) Describe and analyze an algorithm that finds the largest sum of elements in a contiguous subarray  $A[i..j]$ . For example, if the array contains the numbers  $(-6, 12, -7, 0, 14, -7, 5)$ , then the largest sum of any contiguous subarray is  $19 = 12 - 7 + 0 + 14$ .
 

-6	12	-7	0	14	-7	5
----	----	----	---	----	----	---

$\underbrace{\phantom{-6 \hspace{1em} 12 \hspace{1em} -7 \hspace{1em} 0 \hspace{1em} 14}}_{19} \hspace{1em} -7 \hspace{1em} 5$
  - (b) Describe and analyze an algorithm that finds the largest *product* of elements in a contiguous subarray  $A[i..j]$ . You may *not* assume that all the numbers in the input array are integers.
2. This series of exercises asks you to develop efficient algorithms to find optimal *subsequences* of various kinds. A subsequence is anything obtained from a sequence by extracting a subset of elements, but keeping them in the same order; the elements of the subsequence need not be contiguous in the original sequence. For example, the strings C, DAMN, YAIOAI, and DYNAMICPROGRAMMING are all subsequences of the sequence DYNAMICPROGRAMMING.
  - (a) Let  $A[1..m]$  and  $B[1..n]$  be two arbitrary arrays. A *common subsequence* of  $A$  and  $B$  is another sequence that is a subsequence of both  $A$  and  $B$ . Describe an efficient algorithm to compute the length of the *longest* common subsequence of  $A$  and  $B$ .
  - (b) Let  $A[1..m]$  and  $B[1..n]$  be two arbitrary arrays. A *common supersequence* of  $A$  and  $B$  is another sequence that contains both  $A$  and  $B$  as subsequences. Describe an efficient algorithm to compute the length of the *shortest* common supersequence of  $A$  and  $B$ .

- (c) Call a sequence  $X[1..n]$  of numbers *oscillating* if  $X[i] < X[i + 1]$  for all even  $i$ , and  $X[i] > X[i + 1]$  for all odd  $i$ . Describe an efficient algorithm to compute the length of the longest oscillating subsequence of an arbitrary array  $A$  of integers.
- (d) Describe an efficient algorithm to compute the length of the shortest oscillating supersequence of an arbitrary array  $A$  of integers.
- (e) Call a sequence  $X[1..n]$  of numbers *accelerating* if  $2 \cdot X[i] < X[i - 1] + X[i + 1]$  for all  $i$ . Describe an efficient algorithm to compute the length of the longest accelerating subsequence of an arbitrary array  $A$  of integers.
- \*(f) Recall that a sequence  $X[1..n]$  of numbers is *increasing* if  $X[i] < X[i + 1]$  for all  $i$ . Describe an efficient algorithm to compute the length of the *longest common increasing subsequence* of two given arrays of integers. For example,  $\langle 1, 4, 5, 6, 7, 9 \rangle$  is the longest common increasing subsequence of the sequences  $\langle 3, 1, 4, 1, 5, 9, 2, 6, 5, 3, 5, 8, 9, 7, 9, 3 \rangle$  and  $\langle 1, 4, 1, 4, 2, 1, 3, 5, 6, 2, 3, 7, 3, 0, 9, 5 \rangle$ .
3. You and your eight-year-old nephew Elmo decide to play a simple card game. At the beginning of the game, the cards are dealt face up in a long row. Each card is worth a different number of points. After all the cards are dealt, you and Elmo take turns removing either the leftmost or rightmost card from the row, until all the cards are gone. At each turn, you can decide which of the two cards to take. The winner of the game is the player that has collected the most points when the game ends.
- Having never taken an algorithms class, Elmo follows the obvious greedy strategy—when it's his turn, Elmo *always* takes the card with the higher point value. Your task is to find a strategy that will beat Elmo whenever possible. (It might seem mean to beat up on a little kid like this, but Elmo absolutely *hates* it when grown-ups let him win.)
- (a) Prove that you should not also use the greedy strategy. That is, show that there is a game that you can win, but only if you do *not* follow the same greedy strategy as Elmo.
- (b) Describe and analyze an algorithm to determine, given the initial sequence of cards, the maximum number of points that you can collect playing against Elmo.
4. It's almost time to show off your flippin' sweet dancing skills! Tomorrow is the big dance contest you've been training for your entire life, except for that summer you spent with your uncle in Alaska hunting wolverines. You've obtained an advance copy of the the list of  $n$  songs that the judges will play during the contest, in chronological order.
- You know all the songs, all the judges, and your own dancing ability extremely well. For each integer  $k$ , you know that if you dance to the  $k$ th song on the schedule, you will be awarded exactly  $Score[k]$  points, but then you will be physically unable to dance for the next  $Wait[k]$  songs (that is, you cannot dance to songs  $k + 1$  through  $k + Wait[k]$ ). The dancer with the highest total score at the end of the night wins the contest, so you want your total score to be as high as possible.
- Describe and analyze an efficient algorithm to compute the maximum total score you can achieve. The input to your sweet algorithm is the pair of arrays  $Score[1..n]$  and  $Wait[1..n]$ .
5. A palindrome is any string that is exactly the same as its reversal, like I, or DEED, or RACECAR, or AMANAPLANACATACANALPANAMA.

- (a) Describe and analyze an algorithm to find the length of the longest *subsequence* of a given string that is also a palindrome. For example, the longest palindrome subsequence of MAHDYNAMICPROGRAMZLETMESHOWYUTHEM is MHYMRORMYHM, so given that string as input, your algorithm should output the number 11.
- (b) Any string can be decomposed into a sequence of palindromes. For example, the string BUBBASEESABANANA ('Bubba sees a banana.') can be broken into palindromes in the following ways (and many others):

$$\begin{aligned} & \text{BUB} + \text{BASEESAB} + \text{ANANA} \\ & \text{B} + \text{U} + \text{BB} + \text{A} + \text{SEES} + \text{ABA} + \text{NAN} + \text{A} \\ & \text{B} + \text{U} + \text{BB} + \text{A} + \text{SEES} + \text{A} + \text{B} + \text{ANANA} \\ & \text{B} + \text{U} + \text{B} + \text{B} + \text{A} + \text{S} + \text{E} + \text{E} + \text{S} + \text{A} + \text{B} + \text{A} + \text{N} + \text{A} + \text{N} + \text{A} \end{aligned}$$

Describe and analyze an efficient algorithm to find the smallest number of palindromes that make up a given input string. For example, given the input string BUBBASEESABANANA, your algorithm would return the integer 3.

6. Suppose you have a subroutine `QUALITY` that can compute the 'quality' of any given string  $A[1..k]$  in  $O(k)$  time. For example, the quality of a string might be 1 if the string is a Québécois curse word, and 0 otherwise.

Given an arbitrary input string  $T[1..n]$ , we would like to break it into contiguous substrings, such that the total quality of all the substrings is as large as possible. For example, the string SAINTCIBOIREDEDESACRAMENTDECRISS can be decomposed into the substrings SAINT + CIBOIRE + DE + SACRAMENT + DE + CRISSE, of which three (or possibly four) are *sacres*.

Describe an algorithm that breaks a string into substrings of maximum total quality, using the `QUALITY` subroutine.

7. Consider two horizontal lines  $\ell_1$  and  $\ell_2$  in the plane. Suppose we are given (the  $x$ -coordinates of)  $n$  distinct points  $a_1, a_2, \dots, a_n$  on  $\ell_1$  and  $n$  distinct points  $b_1, b_2, \dots, b_n$  on  $\ell_2$ . (The points  $a_i$  and  $b_j$  are *not* necessarily indexed in order from left to right, or in the same order.) Design an algorithm to compute the largest set  $S$  of non-intersecting line segments satisfying to the following restrictions:

- (a) Each segment in  $S$  connects some point  $a_i$  to the corresponding point  $b_i$ .
  - (b) No two segments in  $S$  intersect.
8. Suppose you are given an  $m \times n$  bitmap, represented by an array  $M[1..m, 1..n]$  whose entries are all 0 or 1. A *solid block* is a subarray of the form  $M[i..i', j..j']$  in which every bit is equal to 1. Describe and analyze an efficient algorithm to find a solid block in  $M$  with maximum area.
9. You are driving a bus along a highway, full of rowdy, hyper, thirsty students and a soda fountain machine. Each minute that a student is on your bus, that student drinks one ounce of soda. Your goal is to drop the students off quickly, so that the total amount of soda consumed by all students is as small as possible.

You know how many students will get off of the bus at each exit. Your bus begins somewhere along the highway (probably not at either end) and move s at a constant speed of 37.4 miles per

hour. You must drive the bus along the highway; however, you may drive forward to one exit then backward to an exit in the opposite direction, switching as often as you like. (You can stop the bus, drop off students, and turn around instantaneously.)

Describe an efficient algorithm to drop the students off so that they drink as little soda as possible. Your input consists of the bus route (a list of the exits, together with the travel time between successive exits), the number of students you will drop off at each exit, and the current location of your bus (which you may assume is an exit).

10. In a previous life, you worked as a cashier in the lost Antarctic colony of Nadira, spending the better part of your day giving change to your customers. Because paper is a very rare and valuable resource in Antarctica, cashiers were required by law to use the fewest bills possible whenever they gave change. Thanks to the numerological predilections of one of its founders, the currency of Nadira, called Dream Dollars, was available in the following denominations: \$1, \$4, \$7, \$13, \$28, \$52, \$91, \$365.<sup>9</sup>
  - (a) The greedy change algorithm repeatedly takes the largest bill that does not exceed the target amount. For example, to make \$122 using the greedy algorithm, we first take a \$91 bill, then a \$28 bill, and finally three \$1 bills. Give an example where this greedy algorithm uses more Dream Dollar bills than the minimum possible.
  - (b) Describe and analyze a recursive algorithm that computes, given an integer  $k$ , the minimum number of bills needed to make  $k$  Dream Dollars. (Don't worry about making your algorithm fast; just make sure it's correct.)
  - (c) Describe a dynamic programming algorithm that computes, given an integer  $k$ , the minimum number of bills needed to make  $k$  Dream Dollars. (This one needs to be fast.)

11. What excitement! The Champaign Spinners and the Urbana Dreamweavers have advanced to meet each other in the World Series of Basket-weaving! The World Champions will be decided by a best-of-  $2n - 1$  series of head-to-head weaving matches, and the first to win  $n$  matches will take home the coveted Golden Basket (for example, a best-of-7 series requiring four match wins, but we will keep the generalized case). We know that for any given match there is a constant probability  $p$  that Champaign will win, and a subsequent probability  $q = 1 - p$  that Urbana will win.

Let  $P(i, j)$  be the probability that Champaign will win the series given that they still need  $i$  more victories, whereas Urbana needs  $j$  more victories for the championship.  $P(0, j) = 1$  for any  $j$ , because Champaign needs no more victories to win. Similarly,  $P(i, 0) = 0$  for any  $i$ , as Champaign cannot possibly win if Urbana already has.  $P(0, 0)$  is meaningless. Champaign wins any particular match with probability  $p$  and loses with probability  $q$ , so

$$P(i, j) = p \cdot P(i - 1, j) + q \cdot P(i, j - 1)$$

for any  $i \geq 1$  and  $j \geq 1$ .

Describe and analyze an efficient algorithm that computes the probability that Champaign will win the series (that is, calculate  $P(n, n)$ ), given the parameters  $n$ ,  $p$ , and  $q$  as input.

---

<sup>9</sup>For more details on the history and culture of Nadira, including images of the various denominations of Dream Dollars, see <http://www.dream-dollars.com>.

12. *Vankin's Mile* is a solitaire game played on an  $n \times n$  square grid. The player starts by placing a token on any square of the grid. Then on each turn, the player moves the token either one square to the right or one square down. The game ends when player moves the token off the edge of the board. Each square of the grid has a numerical value, which could be positive, negative, or zero. The player starts with a score of zero; whenever the token lands on a square, the player adds its value to his score. The object of the game is to score as many points as possible.

For example, given the grid below, the player can score  $8 - 6 + 7 - 3 + 4 = 10$  points by placing the initial token on the 8 in the second row, and then moving down, down, right, down, down. (This is *not* the best possible score for these values.)

-1	7	-8	10	-5
-4	-9	<b>8</b>	-6	0
5	-2	<b>-6</b>	-6	7
-7	4	<b>7</b> $\Rightarrow$ <b>-3</b>	-3	-3
7	1	-6	<b>4</b>	-9

Describe and analyze an efficient algorithm to compute the maximum possible score for a game of Vankin's Mile, given the  $n \times n$  array of values as input.

13. A *shuffle* of two strings  $X$  and  $Y$  is formed by interspersing the characters into a new string, keeping the characters of  $X$  and  $Y$  in the same order. For example, ‘bananaananas’ is a shuffle of ‘banana’ and ‘ananas’ in several different ways.

banana      ban      b  
        ananas    ana ana    an a na na s

The strings ‘prodgyrnamatammicng’ and ‘dyprongarmammicing’ are both shuffles of ‘dynamic’ and ‘programming’:

pro<sup>d</sup>g y r<sup>n</sup>a m<sup>m</sup> a m<sup>i</sup>c<sup>g</sup>      dy p r o n g a r<sup>m</sup> a m<sup>m</sup> i c i n g

Given three strings  $A[1..m]$ ,  $B[1..n]$ , and  $C[1..m+n]$ , describe and analyze an algorithm to determine whether  $C$  is a shuffle of  $A$  and  $B$ .

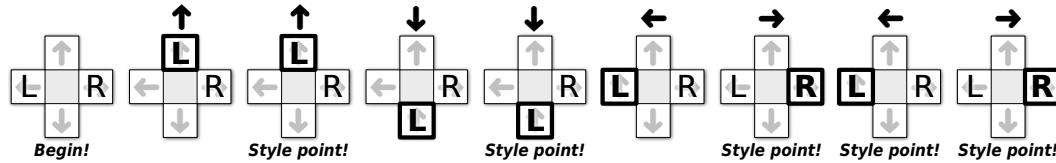
14. Describe and analyze an efficient algorithm to find the length of the longest contiguous substring that appears both forward and backward in an input string  $T[1..n]$ . The forward and backward substrings must not overlap. Here are several examples:

- Given the input string ALGORITHM, your algorithm should return 0.
- Given the input string RECURSION, your algorithm should return 1, for the substring R.
- Given the input string REDIVIDE, your algorithm should return 3, for the substring EDI. (The forward and backward substrings must not overlap!)
- Given the input string DYNAMICPROGRAMMINGMANYTIMES, your algorithm should return 4, for the substring YNAM. (In particular, it should *not* return 6, for the subsequence YNAMIR).

15. *Dance Dance Revolution* is a dance video game, first introduced in Japan by Konami in 1998. Players stand on a platform marked with four arrows, pointing forward, back, left, and right, arranged in a cross pattern. During play, the game plays a song and scrolls a sequence of  $n$  arrows ( $\leftarrow$ ,  $\uparrow$ ,  $\downarrow$ , or  $\rightarrow$ ) from the bottom to the top of the screen. At the precise moment each arrow reaches the top of the screen, the player must step on the corresponding arrow on the dance platform. (The arrows are timed so that you'll step with the beat of the song.)

You are playing a variant of this game called “Vogue Vogue Revolution”, where the goal is to play perfectly but move as little as possible. When an arrow reaches the top of the screen, if one of your feet is already on the correct arrow, you are awarded one style point for maintaining your current pose. If neither foot is on the right arrow, you must move one (and *only* one) of your feet from its current location to the correct arrow on the platform. If you ever step on the wrong arrow, or fail to step on the correct arrow, or move more than one foot at a time, or move either foot when you are already standing on the correct arrow, all your style points are taken away and you lose the game.

How should you move your feet to maximize your total number of style points? For purposes of this problem, assume you always start with your left foot on  $\leftarrow$  and your right foot on  $\rightarrow$ , and that you’ve memorized the entire sequence of arrows. For example, if the sequence is  $\uparrow\uparrow\downarrow\downarrow\leftarrow\rightarrow\leftarrow\rightarrow$ , you can earn 5 style points by moving your feet as shown below:



- (a) **Prove** that for *any* sequence of  $n$  arrows, it is possible to earn at least  $n/4 - 1$  style points.
  - (b) Describe an efficient algorithm to find the maximum number of style points you can earn during a given VVR routine. The input to your algorithm is an array  $\text{Arrow}[1..n]$  containing the sequence of arrows.
16. Suppose we want to display a paragraph of text on a computer screen. The text consists of  $n$  words, where the  $i$ th word is  $p_i$  pixels wide. We want to break the paragraph into several lines, each exactly  $P$  pixels long. Depending on which words we put on each line, we will need to insert different amounts of white space between the words. The paragraph should be fully justified, meaning that the first word on each line starts at its leftmost pixel, and *except for the last line*, the last character on each line ends at its rightmost pixel. There must be at least one pixel of white-space between any two words on the same line.

Define the *slop* of a paragraph layout as the sum over all lines, *except the last*, of the cube of the number of extra white-space pixels in each line (not counting the one pixel required between every adjacent pair of words). Specifically, if a line contains words  $i$  through  $j$ , then the amount of extra white space on that line is  $P - j + i - \sum_{k=i}^j p_k$ . Describe a dynamic programming algorithm to print the paragraph with minimum slop.

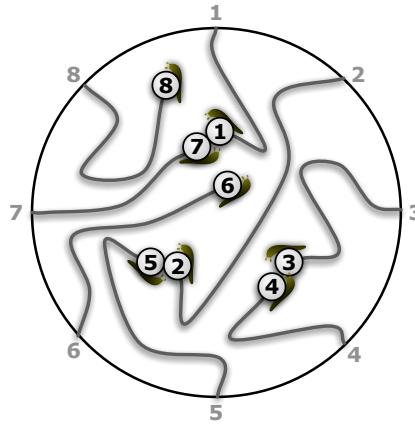
17. [*Stolen from abhi shelat.*] You have mined a large slab of marble from your quarry. For simplicity, suppose the marble slab is a rectangle measuring  $n$  inches in height and  $m$  inches in width. You want to cut the slab into smaller rectangles of various sizes—some for kitchen countertops, some

for large sculpture projects, others for memorial headstones. You have a marble saw that can make either horizontal or vertical cuts across any rectangular slab. At any given time, you can query the spot price  $p_{x,y}$  of an  $x \times y$  marble rectangle, for any positive integers  $x$  and  $y$ . These prices will vary with demand, so do not make any assumptions about them; in particular, larger rectangles may have much smaller spot prices. Given the spot prices, describe an algorithm to compute how to subdivide an  $n \times m$  marble slab to maximize your profit.

18. Every year, as part of its annual meeting, the Antarctic Snail Lovers of Upper Glacierville hold a Round Table Mating Race. Several high-quality breeding snails are placed at the edge of a round table. The snails are numbered in order around the table from 1 to  $n$ . During the race, each snail wanders around the table, leaving a trail of slime behind it. The snails have been specially trained never to fall off the edge of the table or to cross a slime trail, even their own. If two snails meet, they are declared a breeding pair, removed from the table, and whisked away to a romantic hole in the ground to make little baby snails. Note that some snails may never find a mate, even if the race goes on forever.

For every pair of snails, the Antarctic SLUG race organizers have posted a monetary reward, to be paid to the owners if that pair of snails meets during the Mating Race. Specifically, there is a two-dimensional array  $M[1..n, 1..n]$  posted on the wall behind the Round Table, where  $M[i, j] = M[j, i]$  is the reward to be paid if snails  $i$  and  $j$  meet.

Describe and analyze an algorithm to compute the maximum total reward that the organizers could be forced to pay, given the array  $M$  as input.



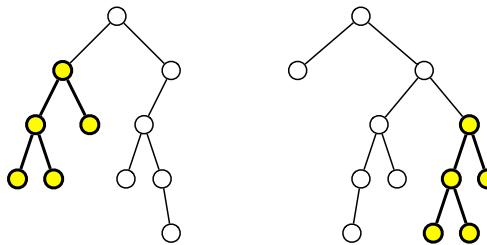
The end of a typical Antarctic SLUG race. Snails 6 and 8 never find mates.

The organizers must pay  $M[3, 4] + M[2, 5] + M[1, 7]$ .

19. Let  $P$  be a set of points in the plane in *convex position*. Intuitively, if a rubber band were wrapped around the points, then every point would touch the rubber band. More formally, for any point  $p$  in  $P$ , there is a line that separates  $p$  from the other points in  $P$ . Moreover, suppose the points are indexed  $P[1], P[2], \dots, P[n]$  in counterclockwise order around the ‘rubber band’, starting with the leftmost point  $P[1]$ .

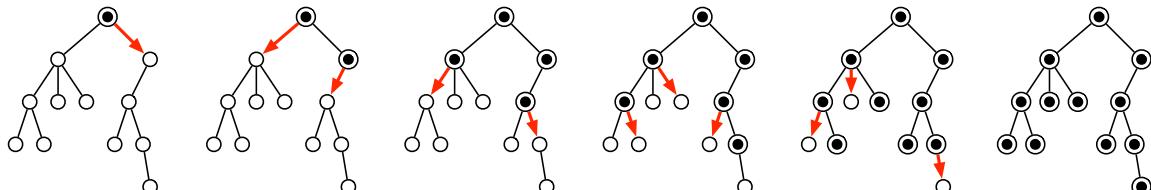
This problem asks you to solve a special case of the traveling salesman problem, where the salesman must visit every point in  $P$ , and the cost of moving from one point  $p \in P$  to another point  $q \in P$  is the Euclidean distance  $|pq|$ .

- (a) Describe a simple algorithm to compute the shortest *cyclic* tour of  $P$ .
- (b) A *simple* tour is one that never crosses itself. Prove that the shortest tour of  $P$  must be simple.
- (c) Describe and analyze an efficient algorithm to compute the shortest tour of  $P$  that starts at the leftmost point  $P[1]$  and ends at the rightmost point  $P[r]$ .
20. Recall that a *subtree* of a (rooted, ordered) binary tree  $T$  consists of a node and all its descendants. Design and analyze an efficient algorithm to compute the **largest common subtree** of two given binary trees  $T_1$  and  $T_2$ ; this is the largest subtree of  $T_1$  that is isomorphic to a subtree in  $T_2$ . The contents of the nodes are irrelevant; we are only interested in matching the underlying combinatorial structure.



Two binary trees, with their largest common subtree emphasized

21. Suppose we need to distribute a message to all the nodes in a rooted tree. Initially, only the root node knows the message. In a single round, any node that knows the message can forward it to at most one of its children. Design an algorithm to compute the minimum number of rounds required for the message to be delivered to all nodes.



A message being distributed through a tree in five rounds.

22. A company is planning a party for its employees. The employees in the company are organized into a strict hierarchy, that is, a tree with the company president at the root. The organizers of the party have assigned a real number to each employee measuring how ‘fun’ the employee is. In order to keep things social, there is one restriction on the guest list: an employee cannot attend the party if their immediate supervisor is present. On the other hand, the president of the company *must* attend the party, even though she has a negative fun rating; it’s her company, after all. Give an algorithm that makes a guest list for the party that maximizes the sum of the ‘fun’ ratings of the guests.

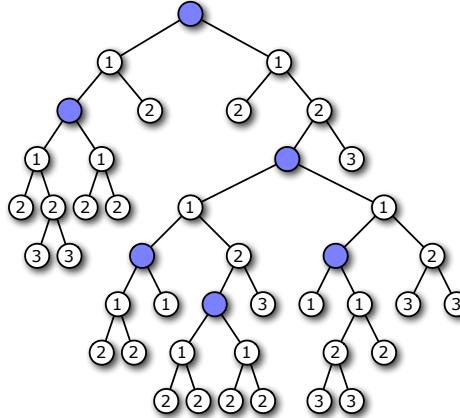
23. Let  $T$  be a rooted binary tree with  $n$  vertices, and let  $k \leq n$  be a positive integer. We would like to mark  $k$  vertices in  $T$  so that every vertex has a nearby marked ancestor. More formally, we define the *clustering cost* of any subset  $K$  of vertices as

$$\text{cost}(K) = \max_v \text{cost}(v, K),$$

where the maximum is taken over all vertices  $v$  in the tree, and

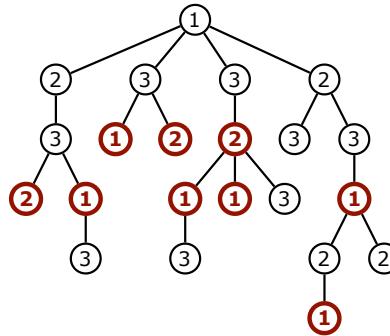
$$\text{cost}(v, K) = \begin{cases} 0 & \text{if } v \in K \\ \infty & \text{if } v \text{ is the root of } T \text{ and } v \notin K \\ 1 + \text{cost}(\text{parent}(v)) & \text{otherwise} \end{cases}$$

Describe and analyze a dynamic-programming algorithm to compute the minimum clustering cost of any subset of  $k$  vertices in  $T$ . For full credit, your algorithm should run in  $O(n^2k^2)$  time.



A subset of 5 vertices with clustering cost 3

24. Oh, no! You have been appointed as the gift czar for Giggle, Inc.'s annual mandatory holiday party! The president of the company, who is certifiably insane, has declared that *every* Giggle employee must receive one of three gifts: (1) an all-expenses-paid six-week vacation anywhere in the world, (2) an all-the-pancakes-you-can-eat breakfast for two at Jumping Jack Flash's Flapjack Stack Shack, or (3) a burning paper bag full of dog poop. Corporate regulations prohibit any employee from receiving the same gift as his/her direct supervisor. Any employee who receives a better gift than his/her direct supervisor will almost certainly be fired in a fit of jealousy. How do you decide what gifts everyone gets if you want to minimize the number of people that get fired?



A tree labeling with cost 9. Bold nodes have smaller labels than their parents.  
This is *not* the optimal labeling for this tree.

More formally, suppose you are given a rooted tree  $T$ , representing the company hierarchy. You want to label each node in  $T$  with an integer 1, 2, or 3, such that every node has a different label from its parent.. The cost of an labeling is the number of nodes that have smaller labels than

their parents. Describe and analyze an algorithm to compute the minimum cost of any labeling of the given tree  $T$ . (Your algorithm does *not* have to compute the actual best labeling—just its cost.)

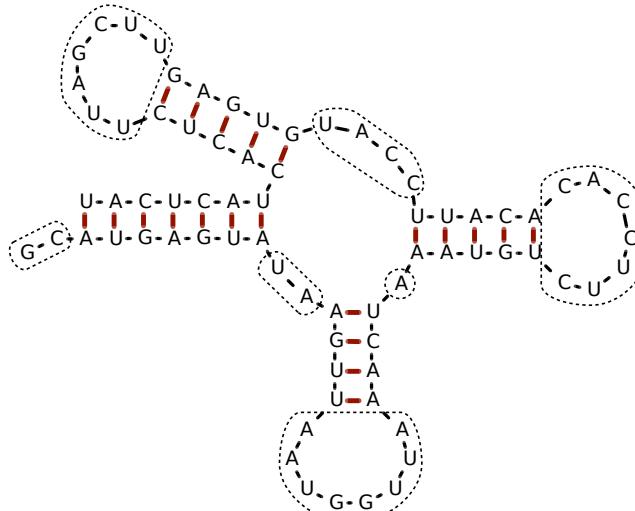
25. Ribonucleic acid (RNA) molecules are long chains of millions of nucleotides or *bases* of four different types: adenine (A), cytosine (C), guanine (G), and uracil (U). The *sequence* of an RNA molecule is a string  $b[1..n]$ , where each character  $b[i] \in \{A, C, G, U\}$  corresponds to a base. In addition to the chemical bonds between adjacent bases in the sequence, hydrogen bonds can form between certain pairs of bases. The set of bonded base pairs is called the *secondary structure* of the RNA molecule.

We say that two base pairs  $(i, j)$  and  $(i', j')$  with  $i < j$  and  $i' < j'$  **overlap** if  $i < i' < j < j'$  or  $i' < i < j' < j$ . In practice, most base pairs are non-overlapping. Overlapping base pairs create so-called *pseudoknots* in the secondary structure, which are essential for some RNA functions, but are more difficult to predict.

Suppose we want to predict the best possible secondary structure for a given RNA sequence. We will adopt a drastically simplified model of secondary structure:

- Each base can be paired with at most one other base.
- Only A-U pairs and C-G pairs can bond.
- Pairs of the form  $(i, i + 1)$  and  $(i, i + 2)$  cannot bond.
- Overlapping base pairs cannot bond.

The last restriction allows us to visualize RNA secondary structure as a sort of fat tree.



Example RNA secondary structure with 21 base pairs, indicated by heavy red lines.  
Gaps are indicated by dotted curves. This structure has score  $2^2 + 2^2 + 8^2 + 1^2 + 7^2 + 4^2 + 7^2 = 187$

- Describe and analyze an algorithm that computes the maximum possible *number* of bonded base pairs in a secondary structure for a given RNA sequence.
- A *gap* in a secondary structure is a maximal substring of unpaired bases. Large gaps lead to chemical instabilities, so secondary structures with smaller gaps are more likely. To account for this preference, let's define the *score* of a secondary structure to be the sum of the *squares*

of the gap lengths.<sup>10</sup> Describe and analyze an algorithm that computes the minimum possible score of a secondary structure for a given RNA sequence.

26. Describe and analyze an algorithm to solve the traveling salesman problem in  $O(2^n \text{poly}(n))$  time. Given an undirected  $n$ -vertex graph  $G$  with weighted edges, your algorithm should return the weight of the lightest cycle in  $G$  that visits every vertex exactly once, or  $\infty$  if  $G$  has no such cycles. [Hint: The obvious recursive algorithm takes  $O(n!)$  time.]

- \*27. Let  $\mathcal{A} = \{A_1, A_2, \dots, A_n\}$  be a finite set of strings over some fixed alphabet  $\Sigma$ . An *edit center* for  $\mathcal{A}$  is a string  $C \in \Sigma^*$  such that the maximum edit distance from  $C$  to any string in  $\mathcal{A}$  is as small as possible. The *edit radius* of  $\mathcal{A}$  is the maximum edit distance from an edit center to a string in  $\mathcal{A}$ . A set of strings may have several edit centers, but its edit radius is unique.

$$\text{EditRadius}(\mathcal{A}) = \min_{C \in \Sigma^*} \max_{A \in \mathcal{A}} \text{Edit}(A, C) \quad \text{EditCenter}(\mathcal{A}) = \arg \min_{C \in \Sigma^*} \max_{A \in \mathcal{A}} \text{Edit}(A, C)$$

- (a) Describe and analyze an efficient algorithm to compute the edit radius of three given strings.
- (b) Describe and analyze an efficient algorithm to approximate the edit radius of an arbitrary set of strings within a factor of 2. (Computing the edit radius exactly is NP-hard unless the number of strings is fixed.)

- \*28. Let  $D[1..n]$  be an array of digits, each an integer between 0 and 9. An *digital subsequence* of  $D$  is a sequence of positive integers composed in the usual way from disjoint substrings of  $D$ . For example, 3, 4, 5, 6, 8, 9, 32, 38, 46, 64, 83, 279 is an increasing digital subsequence of the first several digits of  $\pi$ :

3, 1, 4, 1, 5, 9, 2, 6, 5, 3, 5, 8, 9, 7, 9, 3, 2, 3, 8, 4, 6, 2, 6, 4, 3, 3, 8, 3, 2, 7, 9

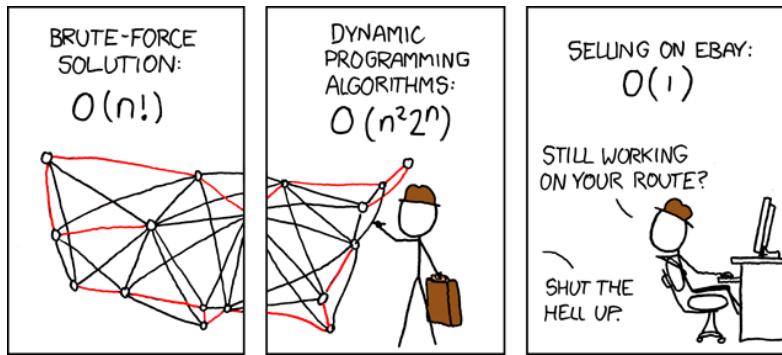
The *length* of a digital subsequence is the number of integers it contains, *not* the number of digits; the preceding example has length 12.

Describe and analyze an efficient algorithm to compute the longest increasing digital subsequence of  $D$ . [Hint: Be careful about your computational assumptions. How long does it take to compare two  $k$ -digit numbers?]

For full credit, your algorithm should run in  $O(n^4)$  time; faster algorithms are worth extra credit. The fastest algorithm I know for this problem runs in  $O(n^{3/2} \log n)$  time, but this requires several tricks, in both the algorithm and its analysis.

---

<sup>10</sup>This score function has absolutely no connection to reality; I just made it up. Real RNA structure prediction requires *much* more complicated scoring functions.



*Ninety percent of science fiction is crud.  
But then, ninety percent of everything is crud,  
and it's the ten percent that isn't crud that is important.*

— [Theodore] Sturgeon's Law (1953)

## \*6 Advanced Dynamic Programming Tricks

Dynamic programming is a powerful technique for efficiently solving recursive problems, but it's hardly the end of the story. In many cases, once we have a basic dynamic programming algorithm in place, we can make further improvements to bring down the running time or the space usage. We saw one example in the Fibonacci number algorithm. Buried inside the naïve iterative Fibonacci algorithm is a recursive problem—computing a power of a matrix—that can be solved more efficiently by dynamic programming techniques—in this case, repeated squaring.

### 6.1 Saving Space: Divide and Conquer

Just as we did for the Fibonacci recurrence, we can reduce the space complexity of our edit distance algorithm from  $O(mn)$  to  $O(m + n)$  by only storing the current and previous rows of the memoization table. This ‘sliding window’ technique provides an easy space improvement for most (but *not all*) dynamic programming algorithm.

Unfortunately, this technique seems to be useful only if we are interested in the *cost* of the optimal edit sequence, not if we want the optimal edit sequence itself. By throwing away most of the table, we apparently lose the ability to walk backward through the table to recover the optimal sequence.

Fortunately for memory-misers, in 1975 Dan Hirshberg discovered a simple divide-and-conquer strategy that allows us to compute the optimal edit sequence in  $O(mn)$  time, using just  $O(m + n)$  space. The trick is to record not just the edit distance for each pair of prefixes, but also a single position in the middle of the optimal editing sequence for that prefix. Specifically, any optimal editing sequence that transforms  $A[1..m]$  into  $B[1..n]$  can be split into two smaller editing sequences, one transforming  $A[1..m/2]$  into  $B[1..h]$  for some integer  $h$ , the other transforming  $A[m/2+1..m]$  into  $B[h+1..n]$ .

To compute this breakpoint  $h$ , we define a second function  $\text{Half}(i, j)$  such that some optimal edit sequence from  $A[1..i]$  into  $B[1..j]$  contains an optimal edit sequence from  $A[1..m/2]$  to  $B[1..\text{Half}(i, j)]$ . We can define this function recursively as follows:

$$\text{Half}(i, j) = \begin{cases} \infty & \text{if } i < m/2 \\ j & \text{if } i = m/2 \\ \text{Half}(i - 1, j) & \text{if } i > m/2 \text{ and } \text{Edit}(i, j) = \text{Edit}(i - 1, j) + 1 \\ \text{Half}(i, j - 1) & \text{if } i > m/2 \text{ and } \text{Edit}(i, j) = \text{Edit}(i, j - 1) + 1 \\ \text{Half}(i - 1, j - 1) & \text{otherwise} \end{cases}$$

(Because there may be more than one optimal edit sequence, this is not the only correct definition.) A simple inductive argument implies that  $\text{Half}(m, n)$  is indeed the correct value of  $h$ . We can easily modify our earlier algorithm so that it computes  $\text{Half}(m, n)$  at the same time as the edit distance  $\text{Edit}(m, n)$ , all in  $O(mn)$  time, using only  $O(m)$  space.

Edit	A	L	G	O	R	I	T	H	M
A	0	1	2	3	4	5	6	7	8
A	1	0	1	2	3	4	5	6	7
L	2	1	0	1	2	3	4	5	6
T	3	2	1	1	2	3	4	4	5
R	4	3	2	2	2	2	3	4	5
U	5	4	3	3	3	3	4	5	6
I	6	5	4	4	4	4	3	4	5
S	7	6	5	5	5	4	4	5	6
T	8	7	6	6	6	5	4	5	6
I	9	8	7	7	7	7	6	5	5
C	10	9	8	8	8	8	7	6	6

Half	A	L	G	O	R	I	T	H	M
A	$\infty$								
A	$\infty$								
L	$\infty$								
T	$\infty$								
R	$\infty$								
U	0	1	2	3	4	5	6	7	8
I	0	1	2	3	4	5	5	5	5
S	0	1	2	3	4	5	5	5	5
T	0	1	2	3	4	5	5	5	5
I	0	1	2	3	4	5	5	5	5
C	0	1	2	3	4	5	5	5	5

Finally, to compute the optimal editing sequence that transforms  $A$  into  $B$ , we recursively compute the optimal sequences transforming  $A[1..m/2]$  into  $B[1..Half(m, n)]$  and transforming  $A[m/2+1..m]$  into  $B[Half(m, n)+1..n]$ . The recursion bottoms out when one string has only constant length, in which case we can determine the optimal editing sequence in linear time using our old dynamic programming algorithm. The running time of the resulting algorithm satisfies the following recurrence:

$$T(m, n) = \begin{cases} O(n) & \text{if } m \leq 1 \\ O(m) & \text{if } n \leq 1 \\ O(mn) + T(m/2, h) + T(m/2, n - h) & \text{otherwise} \end{cases}$$

It's easy to prove inductively that  $T(m, n) = O(mn)$ , no matter what the value of  $h$  is. Specifically, the entire algorithm's running time is at most twice the time for the initial dynamic programming phase.

$$\begin{aligned} T(m, n) &\leq \alpha mn + T(m/2, h) + T(m/2, n - h) \\ &\leq \alpha mn + 2\alpha mh/2 + 2\alpha m(n - h)/2 && [\text{inductive hypothesis}] \\ &= 2\alpha mn \end{aligned}$$

A similar inductive argument implies that the algorithm uses only  $O(n + m)$  space.

Hirschberg's divide-and-conquer trick can be applied to almost any dynamic programming problem to obtain an algorithm to construct an optimal *structure* (in this case, the cheapest edit sequence) within the same space and time bounds as computing the *cost* of that optimal structure (in this case, edit distance). For this reason, we will almost always ask you for algorithms to compute the cost of some optimal structure, not the optimal structure itself.

## 6.2 Saving Time: Sparseness

In many applications of dynamic programming, we are faced with instances where almost every recursive subproblem will be resolved exactly the same way. We call such instances *spARSE*. For example, we might want to compute the edit distance between two strings that have few characters in common, which means there are few “free” substitutions anywhere in the table. Most of the table has exactly the same structure. If we can reconstruct the entire table from just a few key entries, then why compute the entire table?

To better illustrate how to exploit sparseness, let's consider a simplification of the edit distance problem, in which substitutions are not allowed (or equivalently, where a substitution counts as two operations instead of one). Now our goal is to maximize the number of “free” substitutions, or equivalently, to find the *longest common subsequence* of the two input strings.

Fix the two input strings  $A[1..n]$  and  $B[1..m]$ . For any indices  $i$  and  $j$ , let  $LCS(i, j)$  denote the length of the longest common subsequence of the prefixes  $A[1..i]$  and  $B[1..j]$ . This function can be defined recursively as follows:

$$LCS(i, j) = \begin{cases} 0 & \text{if } i = 0 \text{ or } j = 0 \\ LCS(i - 1, j - 1) + 1 & \text{if } A[i] = B[j] \\ \max \{LCS(i, j - 1), LCS(i - 1, j)\} & \text{otherwise} \end{cases}$$

This recursive definition directly translates into an  $O(mn)$ -time dynamic programming algorithm.

Call an index pair  $(i, j)$  a *match point* if  $A[i] = B[j]$ . In some sense, match points are the only ‘interesting’ locations in the memoization table; given a list of the match points, we could easily reconstruct the entire table.

« A L G O R I T H M S »														
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0														
A	0	1	1	1	1	1	1	1	1	1	1	1	1	1
L	0	0	2	2	2	2	2	2	2	2	2	2	2	2
T	0	0	0	2	2	2	2	3	3	3	3	3	3	3
R	0	0	0	0	3	3	3	3	3	3	3	3	3	3
U	0	0	0	0	0	3	3	3	3	3	3	3	3	3
I	0	0	0	0	0	4	4	4	4	4	4	4	4	4
S	0	0	0	0	0	0	4	4	4	4	5	5	5	5
T	0	0	0	0	0	0	0	5	5	5	5	5	5	5
I	0	0	0	0	0	0	0	0	5	5	5	5	5	5
C	0	0	0	0	0	0	0	0	0	5	5	5	5	5
»	0	0	0	0	0	0	0	0	0	0	5	5	5	6

The  $LCS$  memoization table for ALGORITHMS and ALTRUISTIC; the brackets « and » are sentinel characters.

More importantly, we can compute the  $LCS$  function directly from the list of match points using the following recurrence:

$$LCS(i, j) = \begin{cases} 0 & \text{if } i = j = 0 \\ \max \{LCS(i', j') \mid A[i'] = B[j'] \text{ and } i' < i \text{ and } j' < j\} + 1 & \text{if } A[i] = B[j] \\ \max \{LCS(i', j') \mid A[i'] = B[j'] \text{ and } i' \leq i \text{ and } j' \leq j\} & \text{otherwise} \end{cases}$$

(Notice that the inequalities are strict in the second case, but not in the third.) To simplify boundary issues, we add unique sentinel characters  $A[0] = B[0]$  and  $A[m + 1] = B[n + 1]$  to both strings. This ensures that the sets on the right side of the recurrence equation are non-empty, and that we only have to consider match points to compute  $LCS(m, n) = LCS(m + 1, n + 1) - 1$ .

If there are  $K$  match points, we can actually compute them all in  $O(m \log m + n \log n + K)$  time. Sort the characters in each input string, but remembering the original index of each character, and then essentially merge the two sorted arrays, as follows:

```

FINDMATCHES( $A[1..m], B[1..n]$ ):
  for  $i \leftarrow 1$  to  $m$ :  $I[i] \leftarrow i$ 
  for  $j \leftarrow 1$  to  $n$ :  $J[j] \leftarrow j$ 
  sort  $A$  and permute  $I$  to match
  sort  $B$  and permute  $J$  to match

   $i \leftarrow 1; j \leftarrow 1$ 
  while  $i < m$  and  $j < n$ 
    if  $A[i] < B[j]$ 
       $i \leftarrow i + 1$ 
    else if  $A[i] > B[j]$ 
       $j \leftarrow j + 1$ 
    else            $\langle\langle$ Found a match! $\rangle\rangle$ 
       $ii \leftarrow i$ 
      while  $A[ii] = A[i]$ 
         $jj \leftarrow j$ 
        while  $B[jj] = B[j]$ 
          report ( $I[ii], J[jj]$ )
           $jj \leftarrow jj + 1$ 
         $ii \leftarrow i + 1$ 
       $i \leftarrow ii; j \leftarrow jj$ 

```

To efficiently evaluate our modified recurrence, we once again turn to dynamic programming. We consider the match points in lexicographic order—the order they would be encountered in a standard row-major traversal of the  $m \times n$  table—so that when we need to evaluate  $LCS(i, j)$ , all match points  $(i', j')$  with  $i' < i$  and  $j' < j$  have already been evaluated.

```

SPARSELCS( $A[1..m], B[1..n]$ ):
   $Match[1..K] \leftarrow \text{FINDMATCHES}(A, B)$ 
   $Match[K + 1] \leftarrow (m + 1, n + 1)$       $\langle\langle$ Add end sentinel $\rangle\rangle$ 
  Sort  $M$  lexicographically
  for  $k \leftarrow 1$  to  $K$ 
     $(i, j) \leftarrow Match[k]$ 
     $LCS[k] \leftarrow 1$                        $\langle\langle$ From start sentinel $\rangle\rangle$ 
    for  $\ell \leftarrow 1$  to  $k - 1$ 
       $(i', j') \leftarrow Match[\ell]$ 
      if  $i' < i$  and  $j' < j$ 
         $LCS[k] \leftarrow \min\{LCS[k], 1 + LCS[\ell]\}$ 
  return  $LCS[K + 1] - 1$ 

```

The overall running time of this algorithm is  $O(m \log m + n \log n + K^2)$ . So as long as  $K = o(\sqrt{mn})$ , this algorithm is actually faster than naïve dynamic programming.

### 6.3 Saving Time: Monotonicity

The SMAWK matrix-searching algorithm is a better example here; the problem is more general, the algorithm is simple, and the proof is self-contained. Next time!

Recall the optimal binary search tree problem from the previous lecture. Given an array  $F[1..n]$  of access frequencies for  $n$  items, the problem is to compute the binary search tree that minimizes the cost of all accesses. A relatively straightforward dynamic programming algorithm solves this problem in  $O(n^3)$  time.

As for longest common subsequence problem, the algorithm can be improved by exploiting some structure in the memoization table. In this case, however, the relevant structure isn't in the table of costs,

but rather in the table used to reconstruct the actual optimal tree. Let  $OptRoot[i, j]$  denote the index of the root of the optimal search tree for the frequencies  $F[i \dots j]$ ; this is always an integer between  $i$  and  $j$ . Donald Knuth proved the following nice monotonicity property for optimal subtrees: If we move either end of the subarray, the optimal root moves in the same direction or not at all. More formally:

$$OptRoot[i, j - 1] \leq OptRoot[i, j] \leq OptRoot[i + 1, j] \text{ for all } i \text{ and } j.$$

This (nontrivial!) observation leads to the following more efficient algorithm:

```
FASTEROPTIMALSEARCHTREE( $f[1 \dots n]$ ):
  INITF( $f[1 \dots n]$ )
  for  $i \leftarrow 1$  downto  $n$ 
     $OptCost[i, i - 1] \leftarrow 0$ 
     $OptRoot[i, i - 1] \leftarrow i$ 
    for  $d \leftarrow 0$  to  $n$ 
      for  $i \leftarrow 1$  to  $n$ 
        COMPUTECOSTANDROOT( $i, i + d$ )
    return  $OptCost[1, n]$ 
```

```
COMPUTECOSTANDROOT( $i, j$ ):
   $OptCost[i, j] \leftarrow \infty$ 
  for  $r \leftarrow OptRoot[i, j - 1]$  to  $OptRoot[i + 1, j]$ 
     $tmp \leftarrow OptCost[i, r - 1] + OptCost[r + 1, j]$ 
    if  $OptCost[i, j] > tmp$ 
       $OptCost[i, j] \leftarrow tmp$ 
       $OptRoot[i, j] \leftarrow r$ 
   $OptCost[i, j] \leftarrow OptCost[i, j] + F[i, j]$ 
```

It's not hard to see that the loop index  $r$  increases monotonically from 1 to  $n$  during each iteration of the *outermost* for loop of FASTEROPTIMALSEARCHTREE. Consequently, the total cost of all calls to COMPUTECOSTANDROOT is only  $O(n^2)$ .

If we formulate the problem slightly differently, this algorithm can be improved even further. Suppose we require the optimum *external* binary tree, where the keys  $A[1 \dots n]$  are all stored at the leaves, and intermediate pivot values are stored at the internal nodes. An algorithm discovered by Ching Hu and Alan Tucker<sup>1</sup> computes the optimal binary search tree in this setting in only  $O(n \log n)$  time!

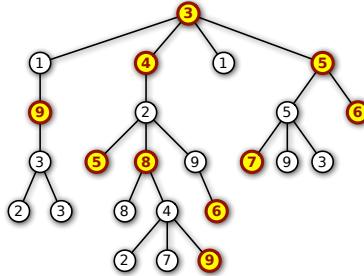
## Exercises

1. Describe an algorithm to compute the edit distance between two strings  $A[1 \dots m]$  and  $B[1 \dots n]$  in  $O(m \log m + n \log n + K^2)$  time, where  $K$  is the number of match points. [Hint: Use the FINDMATCHES algorithm on page 3 as a subroutine.]
  
2. (a) Describe an algorithm to compute the *longest increasing subsequence* of a string  $X[1 \dots n]$  in  $O(n \log n)$  time.  
(b) Using your solution to part (a) as a subroutine, describe an algorithm to compute the longest common subsequence of two strings  $A[1 \dots m]$  and  $B[1 \dots n]$  in  $O(m \log m + n \log n + K \log K)$  time, where  $K$  is the number of match points.
  
3. Describe an algorithm to compute the edit distance between two strings  $A[1 \dots m]$  and  $B[1 \dots n]$  in  $O(m \log m + n \log n + K \log K)$  time, where  $K$  is the number of match points. [Hint: Combine your answers for problems 1 and 2(b).]

<sup>1</sup>T. C. Hu and A. C. Tucker, Optimal computer search trees and variable length alphabetic codes, *SIAM J. Applied Math.* 21:514–532, 1971. For a slightly simpler algorithm with the same running time, see A. M. Garsia and M. L. Wachs, A new algorithms for minimal binary search trees, *SIAM J. Comput.* 6:622–642, 1977. The original correctness proofs for both algorithms are rather intricate; for simpler proofs, see Marek Karpinski, Lawrence L. Larmore, and Wojciech Rytter, Correctness of constructing optimal alphabetic trees revisited, *Theoretical Computer Science*, 180:309–324, 1997.

4. Let  $T$  be an arbitrary rooted tree, where each vertex is labeled with a positive integer. A subset  $S$  of the nodes of  $T$  is *heap-ordered* if it satisfies two properties:

- $S$  contains a node that is an ancestor of every other node in  $S$ .
- For any node  $v$  in  $S$ , the label of  $v$  is larger than the labels of any ancestor of  $v$  in  $S$ .



A heap-ordered subset of nodes in a tree.

- Describe an algorithm to find the largest heap-ordered subset  $S$  of nodes in  $T$  that has the heap property in  $O(n^2)$  time.
- Modify your algorithm from part (a) so that it runs in  $O(n \log n)$  time when  $T$  is a linked list.  
*[Hint: This special case is equivalent to a problem you've seen before.]*
- \*Describe an algorithm to find the largest subset  $S$  of nodes in  $T$  that has the heap property, in  $O(n \log n)$  time. *[Hint: Find an algorithm to merge two sorted lists of lengths  $k$  and  $\ell$  in  $O(\log \binom{k+\ell}{k})$  time.]*

*The point is, ladies and gentleman, greed is good. Greed works, greed is right. Greed clarifies, cuts through, and captures the essence of the evolutionary spirit. Greed in all its forms, greed for life, money, love, knowledge has marked the upward surge in mankind. And greed—mark my words—will save not only Teldar Paper but the other malfunctioning corporation called the USA.*

— Michael Douglas as Gordon Gekko, *Wall Street* (1987)

*There is always an easy solution to every human problem—neat, plausible, and wrong.*

— H. L. Mencken, “The Divine Afflatus”,  
*New York Evening Mail* (November 16, 1917)

## 7 Greedy Algorithms

### 7.1 Storing Files on Tape

Suppose we have a set of  $n$  files that we want to store on a tape. In the future, users will want to read those files from the tape. Reading a file from tape isn’t like reading a file from disk; first we have to fast-forward past all the other files, and that takes a significant amount of time. Let  $L[1..n]$  be an array listing the lengths of each file; specifically, file  $i$  has length  $L[i]$ . If the files are stored in order from 1 to  $n$ , then the cost of accessing the  $k$ th file is

$$\text{cost}(k) = \sum_{i=1}^k L[i].$$

The cost reflects the fact that before we read file  $k$  we must first scan past all the earlier files on the tape. If we assume for the moment that each file is equally likely to be accessed, then the *expected* cost of searching for a random file is

$$\mathbb{E}[\text{cost}] = \sum_{k=1}^n \frac{\text{cost}(k)}{n} = \sum_{k=1}^n \sum_{i=1}^k \frac{L[i]}{n}.$$

If we change the order of the files on the tape, we change the cost of accessing the files; some files become more expensive to read, but others become cheaper. Different file orders are likely to result in different expected costs. Specifically, let  $\pi(i)$  denote the index of the file stored at position  $i$  on the tape. Then the expected cost of the permutation  $\pi$  is

$$\mathbb{E}[\text{cost}(\pi)] = \sum_{k=1}^n \sum_{i=1}^k \frac{L[\pi(i)]}{n}.$$

Which order should we use if we want the expected cost to be as small as possible? The answer is intuitively clear; we should store the files in order from shortest to longest. So let’s prove this.

**Lemma 1.**  $\mathbb{E}[\text{cost}(\pi)]$  is minimized when  $L[\pi(i)] \leq L[\pi(i+1)]$  for all  $i$ .

**Proof:** Suppose  $L[\pi(i)] > L[\pi(i+1)]$  for some  $i$ . To simplify notation, let  $a = \pi(i)$  and  $b = \pi(i+1)$ . If we swap files  $a$  and  $b$ , then the cost of accessing  $a$  increases by  $L[b]$ , and the cost of accessing  $b$  decreases by  $L[a]$ . Overall, the swap changes the expected cost by  $(L[b] - L[a])/n$ . But this change is an improvement, because  $L[b] < L[a]$ . Thus, if the files are out of order, we can improve the expected cost by swapping some mis-ordered adjacent pair.  $\square$

This example gives us our first *greedy algorithm*. To minimize the *total* expected cost of accessing the files, we put the file that is cheapest to access first, and then recursively write everything else; no backtracking, no dynamic programming, just make the best local choice and blindly plow ahead. If we use an efficient sorting algorithm, the running time is clearly  $O(n \log n)$ , plus the time required to actually write the files. To prove the greedy algorithm is actually correct, we simply prove that the output of any other algorithm can be improved by some sort of swap.

Let's generalize this idea further. Suppose we are also given an array  $f[1..n]$  of *access frequencies* for each file; file  $i$  will be accessed exactly  $f[i]$  times over the lifetime of the tape. Now the *total* cost of accessing all the files on the tape is

$$\Sigma\text{cost}(\pi) = \sum_{k=1}^n \left( f[\pi(k)] \cdot \sum_{i=1}^k L[\pi(i)] \right) = \sum_{k=1}^n \sum_{i=1}^k (f[\pi(k)] \cdot L[\pi(i)]).$$

Now what order should store the files if we want to minimize the total cost?

We've already proved that if all the frequencies are equal, then we should sort the files by increasing size. If the frequencies are all different but the file lengths  $L[i]$  are all equal, then intuitively, we should sort the files by *decreasing* access frequency, with the most-accessed file first. In fact, this is not hard to prove by modifying the proof of Lemma 1. But what if the sizes and the frequencies are both different? In this case, we should sort the files by the ratio  $L/f$ .

**Lemma 2.**  $\Sigma\text{cost}(\pi)$  is minimized when  $\frac{L[\pi(i)]}{F[\pi(i)]} \leq \frac{L[\pi(i+1)]}{F[\pi(i+1)]}$  for all  $i$ .

**Proof:** Suppose  $L[\pi(i)]/F[\pi(i)] > L[\pi(i+1)]/F[\pi(i+1)]$  for some  $i$ . To simplify notation, let  $a = \pi(i)$  and  $b = \pi(i+1)$ . If we swap files  $a$  and  $b$ , then the cost of accessing  $a$  increases by  $L[b]$ , and the cost of accessing  $b$  decreases by  $L[a]$ . Overall, the swap changes the total cost by  $L[b]F[a] - L[a]F[a]$ . But this change is an improvement, since

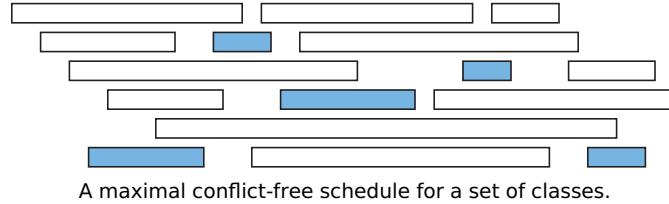
$$\frac{L[a]}{F[a]} > \frac{L[b]}{F[b]} \implies L[b]F[a] - L[a]F[a] < 0.$$

Thus, if the files are out of order, we can improve the total cost by swapping some mis-ordered adjacent pair.  $\square$

## 7.2 Scheduling Classes

The next example is slightly less trivial. Suppose you decide to drop out of computer science at the last minute and change your major to Applied Chaos. The Applied Chaos department offers all of its classes on the same day every week, called ‘Soberday’ by the students.<sup>1</sup> Every class has a different start time and a different ending time: AC 101 (‘Toilet Paper Landscape Architecture’) starts at 10:27pm and ends at 11:51pm; AC 666 (‘Immanentizing the Eschaton’) starts at 4:18pm and ends at 7:06pm, and so on. In the interest of graduating as quickly as possible, you want to register for as many classes as you can. (Applied Chaos classes don’t require any actual work.) The university’s registration computer won’t let you register for overlapping classes, and no one in the department knows how to override this ‘feature’. Which classes should you take?

More formally, suppose you are given two arrays  $S[1..n]$  and  $F[1..n]$  listing the start and finish times of each class. Your task is to choose the largest possible subset  $X \in \{1, 2, \dots, n\}$  so that for any pair  $i, j \in X$ , either  $S[i] > F[j]$  or  $S[j] > F[i]$ . We can illustrate the problem by drawing each class as a



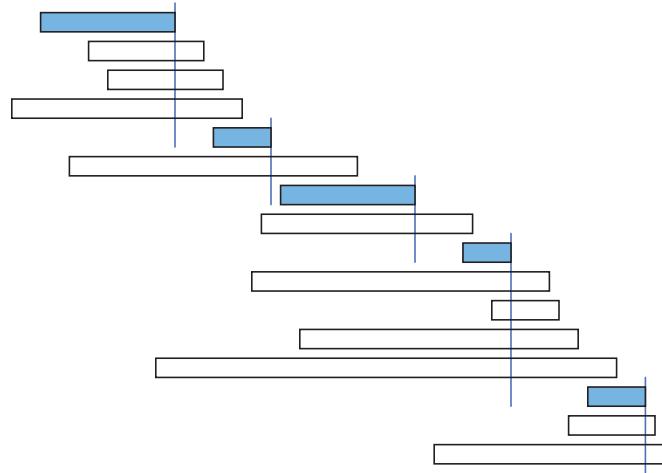
rectangle whose left and right  $x$ -coordinates show the start and finish times. The goal is to find a largest subset of rectangles that do not overlap vertically.

This problem has a fairly simple recursive solution, based on the observation that either you take class 1 or you don't. Let  $B_4$  denote the set of classes that end *before* class 1 starts, and let  $L_8$  denote the set of classes that start *later* than class 1 ends:

$$B_4 = \{i \mid 2 \leq i \leq n \text{ and } F[i] < S[1]\} \quad L_8 = \{i \mid 2 \leq i \leq n \text{ and } S[i] > F[1]\}$$

If class 1 is in the optimal schedule, then so are the optimal schedules for  $B_4$  and  $L_8$ , which we can find recursively. If not, we can find the optimal schedule for  $\{2, 3, \dots, n\}$  recursively. So we should try both choices and take whichever one gives the better schedule. Evaluating this recursive algorithm from the bottom up gives us a dynamic programming algorithm that runs in  $O(n^2)$  time. I won't bother to go through the details, because we can do better.<sup>2</sup>

Intuitively, we'd like the first class to finish as early as possible, because that leaves us with the most remaining classes. If this greedy strategy works, it suggests the following very simple algorithm. Scan through the classes in order of finish time; whenever you encounter a class that doesn't conflict with your latest class so far, take it!



The same classes sorted by finish times and the greedy schedule.

We can write the greedy algorithm somewhat more formally as follows. (Hopefully the first line is understandable.)

---

<sup>1</sup>but interestingly, *not* by the faculty

<sup>2</sup>But you should still work out the details yourself. The dynamic programming algorithm can be used to find the "best" schedule for several different definitions of "best", but the greedy algorithm I'm about to describe only works when "best" means "biggest". Also, you need the practice.

```

GREEDYSCHEDULE( $S[1..n], F[1..n]$ ):
    sort  $F$  and permute  $S$  to match
     $count \leftarrow 1$ 
     $X[count] \leftarrow 1$ 
    for  $i \leftarrow 2$  to  $n$ 
        if  $S[i] > F[X[count]]$ 
             $count \leftarrow count + 1$ 
             $X[count] \leftarrow i$ 
    return  $X[1..count]$ 

```

This algorithm clearly runs in  $O(n \log n)$  time.

To prove that this algorithm actually gives us a maximal conflict-free schedule, we use an exchange argument, similar to the one we used for tape sorting. We are not claiming that the greedy schedule is the *only* maximal schedule; there could be others. (See the figures on the previous page.) All we can claim is that at least one of the maximal schedules is the one that the greedy algorithm produces.

**Lemma 3.** *At least one maximal conflict-free schedule includes the class that finishes first.*

**Proof:** Let  $f$  be the class that finishes first. Suppose we have a maximal conflict-free schedule  $X$  that does not include  $f$ . Let  $g$  be the first class in  $X$  to finish. Since  $f$  finishes before  $g$  does,  $f$  cannot conflict with any class in the set  $S \setminus \{g\}$ . Thus, the schedule  $X' = X \cup \{f\} \setminus \{g\}$  is also conflict-free. Since  $X'$  has the same size as  $X$ , it is also maximal.  $\square$

To finish the proof, we call on our old friend, induction.

**Theorem 4.** *The greedy schedule is an optimal schedule.*

**Proof:** Let  $f$  be the class that finishes first, and let  $L$  be the subset of classes the start after  $f$  finishes. The previous lemma implies that some optimal schedule contains  $f$ , so the best schedule that contains  $f$  is an optimal schedule. The best schedule that includes  $f$  must contain an optimal schedule for the classes that do not conflict with  $f$ , that is, an optimal schedule for  $L$ . The greedy algorithm chooses  $f$  and then, by the inductive hypothesis, computes an optimal schedule of classes from  $L$ .  $\square$

The proof might be easier to understand if we unroll the induction slightly.

**Proof:** Let  $\langle g_1, g_2, \dots, g_k \rangle$  be the sequence of classes chosen by the greedy algorithm. Suppose we have a maximal conflict-free schedule of the form

$$\langle g_1, g_2, \dots, g_{j-1}, c_j, c_{j+1}, \dots, c_m \rangle,$$

where the classes  $c_i$  are different from the classes chosen by the greedy algorithm. By construction, the  $j$ th greedy choice  $g_j$  does not conflict with any earlier class  $g_1, g_2, \dots, g_{j-1}$ , and since our schedule is conflict-free, neither does  $c_j$ . Moreover,  $g_j$  has the *earliest* finish time among all classes that don't conflict with the earlier classes; in particular,  $g_j$  finishes before  $c_j$ . This implies that  $g_j$  does not conflict with any of the later classes  $c_{j+1}, \dots, c_m$ . Thus, the schedule

$$\langle g_1, g_2, \dots, g_{j-1}, g_j, c_{j+1}, \dots, c_m \rangle,$$

is conflict-free. (This is just a generalization of Lemma 3, which considers the case  $j = 1$ .) By induction, it now follows that there is an optimal schedule  $\langle g_1, g_2, \dots, g_k, c_{k+1}, \dots, c_m \rangle$  that includes every class chosen by the greedy algorithm. But this is impossible unless  $k = m$ ; if there were a class  $c_{k+1}$  that does not conflict with  $g_k$ , the greedy algorithm would choose more than  $k$  classes.  $\square$

### 7.3 General Structure

The basic structure of this correctness proof is exactly the same as for the tape-sorting problem: an inductive exchange argument.

- Assume that there is an optimal solution that is different from the greedy solution.
- Find the ‘first’ difference between the two solutions.
- Argue that we can exchange the optimal choice for the greedy choice without degrading the solution.

This argument implies by induction that there is an optimal solution that contains the entire greedy solution. Sometimes, as in the scheduling problem, an additional step is required to show no optimal solution *strictly* improves the greedy solution.

### 7.4 Huffman Codes

A *binary code* assigns a string of 0s and 1s to each character in the alphabet. A binary code is *prefix-free* if no code is a prefix of any other. 7-bit ASCII and Unicode’s UTF-8 are both prefix-free binary codes. Morse code is a binary code, but it is not prefix-free; for example, the code for S (· · ·) includes the code for E (·) as a prefix. Any prefix-free binary code can be visualized as a binary tree with the encoded characters stored at the leaves. The code word for any symbol is given by the path from the root to the corresponding leaf; 0 for left, 1 for right. The length of a codeword for a symbol is the depth of the corresponding leaf. (Note that the code tree is *not* a binary search tree. We don’t care at all about the sorted order of symbols at the leaves. (In fact, the symbols may not have a well-defined order!)

Suppose we want to encode messages in an  $n$ -character alphabet so that the encoded message is as short as possible. Specifically, given an array frequency counts  $f[1..n]$ , we want to compute a prefix-free binary code that minimizes the total encoded length of the message:<sup>3</sup>

$$\sum_{i=1}^n f[i] \cdot \text{depth}(i).$$

In 1951, as a PhD student at MIT, David Huffman developed the following greedy algorithm to produce such an optimal code:<sup>4</sup>

HUFFMAN: Merge the two least frequent letters and recurse.

For example, suppose we want to encode the following helpfully self-descriptive sentence, discovered by Lee Sallows.<sup>5</sup>

This sentence contains three a’s, three c’s, two d’s, twenty-six e’s, five f’s, three g’s, eight h’s, thirteen i’s, two l’s, sixteen n’s, nine o’s, six r’s, twenty-seven s’s, twenty-two t’s, two u’s, five v’s, eight w’s, four x’s, five y’s, and only one z.

---

<sup>3</sup>This looks almost exactly like the cost of a binary search tree, but the optimization problem is very different: code trees are not required to keep the keys in any particular order.

<sup>4</sup>Huffman was a student in an information theory class taught by Robert Fano, who was a close colleague of Claude Shannon, the father of information theory. Fano and Shannon had previously developed a different greedy algorithm for producing prefix codes—split the frequency array into two subarrays as evenly as possible, and then recursively builds a code for each subarray—but these Fano-Shannon codes were known not to be optimal. Fano posed the (then open) problem of finding an optimal encoding to his class; Huffman solved the problem as a class project, in lieu of taking a final exam.

<sup>5</sup>A. K. Dewdney. Computer recreations. *Scientific American*, October 1984. Douglas Hofstadter published a few earlier examples of Lee Sallows’ self-descriptive sentences in his *Scientific American* column in January 1982.

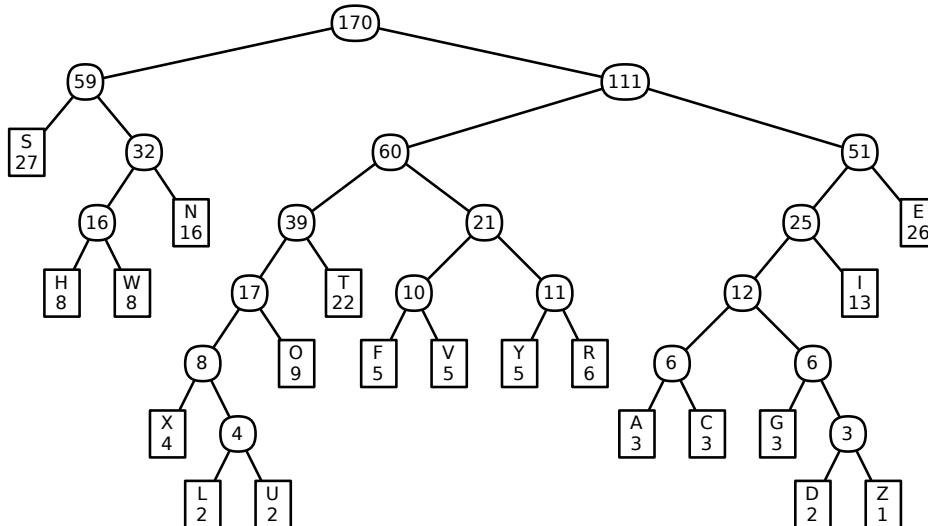
To keep things simple, let's forget about the forty-four spaces, nineteen apostrophes, nineteen commas, three hyphens, and only one period, and just encode the letters. Here's the frequency table:

A	C	D	E	F	G	H	I	L	N	O	R	S	T	U	V	W	X	Y	Z
3	3	2	26	5	3	8	13	2	16	9	6	27	22	2	5	8	4	5	1

Huffman's algorithm picks out the two least frequent letters, breaking ties arbitrarily—in this case, say, Z and D—and merges them together into a single new character  $\Delta$  with frequency 3. This new character becomes an internal node in the code tree we are constructing, with Z and D as its children; it doesn't matter which child is which. The algorithm then recursively constructs a Huffman code for the new frequency table

A	C	E	F	G	H	I	L	N	O	R	S	T	U	V	W	X	Y	$\Delta$
3	3	26	5	3	8	13	2	16	9	6	27	22	2	5	8	4	5	3

After 19 merges, all 20 characters have been merged together. The record of merges gives us our code tree. The algorithm makes a number of arbitrary choices; as a result, there are actually several different Huffman codes. One such code is shown below. For example, the code for A is 110000, and the code for S is 00.



A Huffman code for Lee Sallows' self-descriptive sentence; the numbers are frequencies for merged characters

If we use this code, the encoded message starts like this:

1001 0100 1101 00 00 111 011 1001 111 011 110001 111 110001 10001 011 1001 110000 1101 ...  
 T H I S S E N T E N C E C O N T A I ...

Here is the list of costs for encoding each character in the example message, along with that character's contribution to the total length of the encoded message:

char.	A	C	D	E	F	G	H	I	L	N	O	R	S	T	U	V	W	X	Y	Z
freq.	3	3	2	26	5	3	8	13	2	16	9	6	27	22	2	5	8	4	5	1
depth	6	6	7	3	5	6	4	4	7	3	4	4	2	4	7	5	4	6	5	7
total	18	18	14	78	25	18	32	52	14	48	36	24	54	88	14	25	32	24	25	7

Altogether, the encoded message is 646 bits long. Different Huffman codes would assign different codes, possibly with different lengths, to various characters, but the overall length of the encoded message is the same for any Huffman code: 646 bits.

Given the simple structure of Huffman's algorithm, it's rather surprising that it produces an *optimal* prefix-free binary code. Encoding Lee Sallows' sentence using *any* prefix-free code requires at least 646 bits! Fortunately, the recursive structure makes this claim easy to prove using an exchange argument, similar to our earlier optimality proofs. We start by proving that the algorithm's very first choice is correct.

**Lemma 5.** *Let  $x$  and  $y$  be the two least frequent characters (breaking ties between equally frequent characters arbitrarily). There is an optimal code tree in which  $x$  and  $y$  are siblings.*

**Proof:** I'll actually prove a stronger statement: There is an optimal code in which  $x$  and  $y$  are siblings and have the largest depth of any leaf.

Let  $T$  be an optimal code tree, and suppose this tree has depth  $d$ . Since  $T$  is a full binary tree, it has at least two leaves at depth  $d$  that are siblings. (Verify this by induction!) Suppose those two leaves are *not*  $x$  and  $y$ , but some other characters  $a$  and  $b$ .

Let  $T'$  be the code tree obtained by swapping  $x$  and  $a$ . The depth of  $x$  increases by some amount  $\Delta$ , and the depth of  $a$  decreases by the same amount. Thus,

$$\text{cost}(T') = \text{cost}(T) - (f[a] - f[x])\Delta.$$

By assumption,  $x$  is one of the two least frequent characters, but  $a$  is not, which implies that  $f[a] \geq f[x]$ . Thus, swapping  $x$  and  $a$  does not increase the total cost of the code. Since  $T$  was an optimal code tree, swapping  $x$  and  $a$  does not decrease the cost, either. Thus,  $T'$  is also an optimal code tree (and incidentally,  $f[a]$  actually equals  $f[x]$ ).

Similarly, swapping  $y$  and  $b$  must give yet another optimal code tree. In this final optimal code tree,  $x$  and  $y$  as maximum-depth siblings, as required.  $\square$

Now optimality is guaranteed by our dear friend the Recursion Fairy! Essentially we're relying on the following recursive definition for a full binary tree: either a single node, or a full binary tree where some leaf has been replaced by an internal node with two leaf children.

**Theorem 6.** *Huffman codes are optimal prefix-free binary codes.*

**Proof:** If the message has only one or two different characters, the theorem is trivial.

Otherwise, let  $f[1..n]$  be the original input frequencies, where without loss of generality,  $f[1]$  and  $f[2]$  are the two smallest. To keep things simple, let  $f[n+1] = f[1] + f[2]$ . By the previous lemma, we know that some optimal code for  $f[1..n]$  has characters 1 and 2 as siblings.

Let  $T'$  be the Huffman code tree for  $f[3..n+1]$ ; the inductive hypothesis implies that  $T'$  is an optimal code tree for the smaller set of frequencies. To obtain the final code tree  $T$ , we replace the leaf labeled  $n+1$  with an internal node with two children, labelled 1 and 2. I claim that  $T$  is optimal for the original frequency array  $f[1..n]$ .

To prove this claim, we can express the cost of  $T$  in terms of the cost of  $T'$  as follows. (In these equations,  $\text{depth}(i)$  denotes the depth of the leaf labelled  $i$  in either  $T$  or  $T'$ ; if the leaf appears in both

$T$  and  $T'$ , it has the same depth in both trees.)

$$\begin{aligned}
 \text{cost}(T) &= \sum_{i=1}^n f[i] \cdot \text{depth}(i) \\
 &= \sum_{i=3}^{n+1} f[i] \cdot \text{depth}(i) + f[1] \cdot \text{depth}(1) + f[2] \cdot \text{depth}(2) - f[n+1] \cdot \text{depth}(n+1) \\
 &= \text{cost}(T') + f[1] \cdot \text{depth}(1) + f[2] \cdot \text{depth}(2) - f[n+1] \cdot \text{depth}(n+1) \\
 &= \text{cost}(T') + (f[1] + f[2]) \cdot \text{depth}(T) - f[n+1] \cdot (\text{depth}(T) - 1) \\
 &= \text{cost}(T') + f[1] + f[2]
 \end{aligned}$$

This equation implies that minimizing the cost of  $T$  is equivalent to minimizing the cost of  $T'$ ; in particular, attaching leaves labeled 1 and 2 to the leaf in  $T'$  labeled  $n+1$  gives an optimal code tree for the original frequencies.  $\square$

To actually implement Huffman codes efficiently, we keep the characters in a min-heap, where the priority of each character is its frequency. We can construct the code tree by keeping three arrays of indices, listing the left and right children and the parent of each node. The root of the tree is the node with index  $2n-1$ .

```

BUILDHUFFMAN( $f[1..n]$ ):
  for  $i \leftarrow 1$  to  $n$ 
     $L[i] \leftarrow 0$ ;  $R[i] \leftarrow 0$ 
    INSERT( $i, f[i]$ )
  for  $i \leftarrow n$  to  $2n-1$ 
     $x \leftarrow \text{EXTRACTMIN}()$ 
     $y \leftarrow \text{EXTRACTMIN}()$ 
     $f[i] \leftarrow f[x] + f[y]$ 
     $L[i] \leftarrow x$ ;  $R[i] \leftarrow y$ 
     $P[x] \leftarrow i$ ;  $P[y] \leftarrow i$ 
    INSERT( $i, f[i]$ )
   $P[2n-1] \leftarrow 0$ 

```

The algorithm performs  $O(n)$  min-heap operations. If we use a balanced binary tree as the heap, each operation requires  $O(\log n)$  time, so the total running time of BUILDHUFFMAN is  $O(n \log n)$ .

Finally, here are simple algorithms to encode and decode messages:

```

HUFFMANENCODE( $A[1..k]$ ):
   $m \leftarrow 1$ 
  for  $i \leftarrow 1$  to  $k$ 
    HUFFMANENCODEONE( $A[i]$ )
HUFFMANENCODEONE( $x$ ):
  if  $x < 2n-1$ 
    HUFFMANENCODEONE( $P[x]$ )
    if  $x = L[P[x]]$ 
       $B[m] \leftarrow 0$ 
    else
       $B[m] \leftarrow 1$ 
     $m \leftarrow m + 1$ 

```

```

HUFFMANDECODE( $B[1..m]$ ):
   $k \leftarrow 1$ 
   $v \leftarrow 2n-1$ 
  for  $i \leftarrow 1$  to  $m$ 
    if  $B[i] = 0$ 
       $v \leftarrow L[v]$ 
    else
       $v \leftarrow R[v]$ 
    if  $L[v] = 0$ 
       $A[k] \leftarrow v$ 
       $k \leftarrow k + 1$ 
     $v \leftarrow 2n-1$ 

```

## Exercises

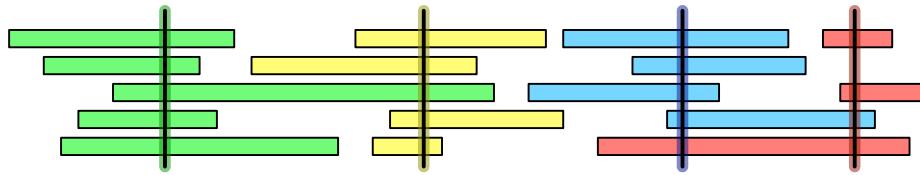
1. Consider the following alternative greedy algorithms for the class scheduling problem. For each algorithm, either prove or disprove that it always constructs an optimal schedule. Assume that all four algorithms break ties arbitrarily.
  - (a) Choose the course that *ends latest*, discard all conflicting classes, and recurse.
  - (b) Choose the course that *starts first*, discard all conflicting classes, and recurse.
  - (c) Choose the course with *shortest duration*, discard all conflicting classes, and recurse.
  - (d) Choose a course that *conflicts with the fewest other courses*, discard all conflicting classes, and recurse.
2. Now consider a weighted version of the class scheduling problem, where different classes offer different number of credit hours (totally unrelated to the duration of the class lectures). Your goal is now to choose a set of non-conflicting classes that give you the largest possible number of credit hours, given an array of start times, end times, and credit hours as input.
  - (a) Prove that the greedy algorithm described in the notes — Choose the class that ends first and recurse — does *not* always return an optimal schedule.
  - (b) Describe an algorithm to compute the optimal schedule in  $O(n^2)$  time.
3. Let  $X$  be a set of  $n$  intervals on the real line. A subset of intervals  $Y \subseteq X$  is called a *tiling path* if the intervals in  $Y$  cover the intervals in  $X$ , that is, any real value that is contained in some interval in  $X$  is also contained in some interval in  $Y$ . The *size* of a tiling cover is just the number of intervals.

Describe and analyze an algorithm to compute the smallest tiling path of  $X$  as quickly as possible. Assume that your input consists of two arrays  $X_L[1..n]$  and  $X_R[1..n]$ , representing the left and right endpoints of the intervals in  $X$ . If you use a greedy algorithm, you must prove that it is correct.



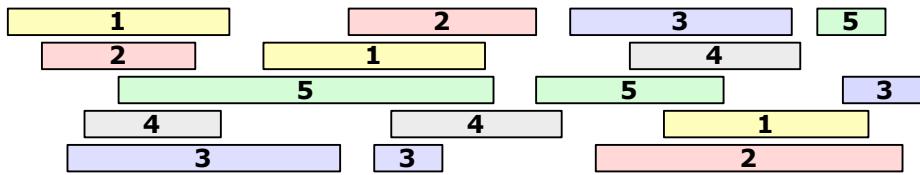
A set of intervals. The seven shaded intervals form a tiling path.

4. Let  $X$  be a set of  $n$  intervals on the real line. We say that a set  $P$  of points *stabs*  $X$  if every interval in  $X$  contains at least one point in  $P$ . Describe and analyze an efficient algorithm to compute the smallest set of points that stabs  $X$ . Assume that your input consists of two arrays  $X_L[1..n]$  and  $X_R[1..n]$ , representing the left and right endpoints of the intervals in  $X$ . As usual, If you use a greedy algorithm, you must prove that it is correct.
5. Let  $X$  be a set of  $n$  intervals on the real line. A *proper coloring* of  $X$  assigns a color to each interval, so that any two overlapping intervals are assigned different colors. Describe and analyze an efficient algorithm to compute the minimum number of colors needed to properly color  $X$ . Assume



A set of intervals stabbed by four points (shown here as vertical segments)

that your input consists of two arrays  $L[1..n]$  and  $R[1..n]$ , where  $L[i]$  and  $R[i]$  are the left and right endpoints of the  $i$ th interval. As usual, if you use a greedy algorithm, you must prove that it is correct.



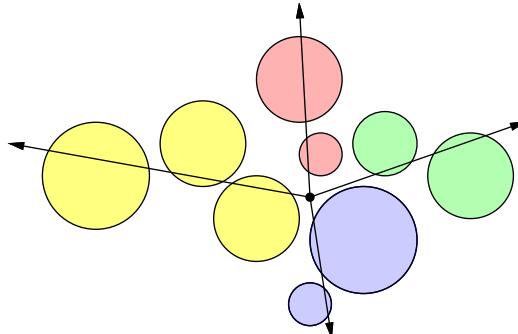
A proper coloring of a set of intervals using five colors.

6. Suppose you are a simple shopkeeper living in a country with  $n$  different types of coins, with values  $1 = c[1] < c[2] < \dots < c[n]$ . (In the U.S., for example,  $n = 6$  and the values are 1, 5, 10, 25, 50 and 100 cents.) Your beloved and benevolent dictator, El Generalissimo, has decreed that whenever you give a customer change, you must use the smallest possible number of coins, so as not to wear out the image of El Generalissimo lovingly engraved on each coin by servants of the Royal Treasury.
  - (a) In the United States, there is a simple greedy algorithm that always results in the smallest number of coins: subtract the largest coin and recursively give change for the remainder. El Generalissimo does not approve of American capitalist greed. Show that there is a set of coin values for which the greedy algorithm does *not* always give the smallest possible number of coins.
  - (b) Now suppose El Generalissimo decides to impose a currency system where the coin denominations are consecutive powers  $b^0, b^1, b^2, \dots, b^k$  of some integer  $b \geq 2$ . Prove that despite El Generalissimo's disapproval, the greedy algorithm described in part (a) does make optimal change in this currency system.
  - (c) Describe and analyze an efficient algorithm to determine, given a target amount  $A$  and a sorted array  $c[1..n]$  of coin denominations, the smallest number of coins needed to make  $A$  cents in change. Assume that  $c[1] = 1$ , so that it is possible to make change for any amount  $A$ .
7. Suppose you have just purchased a new type of hybrid car that uses fuel extremely efficiently, but can only travel 100 miles on a single battery. The car's fuel is stored in a single-use battery, which must be replaced after at most 100 miles. The actual fuel is virtually free, but the batteries are expensive and can only be installed by licensed battery-replacement technicians. Thus, even if you decide to replace your battery early, you must still pay full price for the new battery to be installed. Moreover, because these batteries are in high demand, no one can afford to own more than one battery at a time.

Suppose you are trying to get from San Francisco to New York City on the new Inter-Continental Super-Highway, which runs in a direct line between these two cities. There are several fueling stations along the way; each station charges a different price for installing a new battery. Before you start your trip, you carefully print the Wikipedia page listing the locations and prices of every fueling station on the ICSH. Given this information, how do you decide the best places to stop for fuel?

More formally, suppose you are given two arrays  $D[1..n]$  and  $C[1..n]$ , where  $D[i]$  is the distance from the start of the highway to the  $i$ th station, and  $C[i]$  is the cost to replace your battery at the  $i$ th station. Assume that your trip starts and ends at fueling stations (so  $D[1] = 0$  and  $D[n]$  is the total length of your trip), and that your car starts with an empty battery (so you must install a new battery at station 1).

- (a) Describe and analyze a greedy algorithm to find the minimum number of refueling stops needed to complete your trip. Don't forget to prove that your algorithm is correct.
  - (b) But what you really want to minimize is the total *cost* of travel. Show that your greedy algorithm in part (a) does *not* produce an optimal solution when extended to this setting.
  - (c) Describe an efficient algorithm to compute the locations of the fuel stations you should stop at to minimize the total cost of travel.
8. Suppose you are standing in a field surrounded by several large balloons. You want to use your brand new Acme Brand Zap-O-Matic™ to pop all the balloons, without moving from your current location. The Zap-O-Matic™ shoots a high-powered laser beam, which pops all the balloons it hits. Since each shot requires enough energy to power a small country for a year, you want to fire as few shots as possible.



Nine balloons popped by 4 shots of the Zap-O-Matic™

The *minimum zap* problem can be stated more formally as follows. Given a set  $C$  of  $n$  circles in the plane, each specified by its radius and the  $(x, y)$  coordinates of its center, compute the minimum number of rays from the origin that intersect every circle in  $C$ . Your goal is to find an efficient algorithm for this problem.

- (a) Suppose it is possible to shoot a ray that does not intersect any balloons. Describe and analyze a greedy algorithm that solves the minimum zap problem in this special case. [Hint: See Exercise 2.]
- (b) Describe and analyze a greedy algorithm whose output is within 1 of optimal. That is, if  $m$  is the minimum number of rays required to hit every balloon, then your greedy algorithm must output either  $m$  or  $m + 1$ . (Of course, you must prove this fact.)

- (c) Describe an algorithm that solves the minimum zap problem in  $O(n^2)$  time.
- \*(d) Describe an algorithm that solves the minimum zap problem in  $O(n \log n)$  time.

Assume you have a subroutine  $\text{INTERSECTS}(r, c)$  that determines whether a ray  $r$  intersects a circle  $c$  in  $O(1)$  time. This subroutine is not difficult to write, but it's not the interesting part of the problem.

*The problem is that we attempt to solve the simplest questions cleverly,  
thereby rendering them unusually complex.  
One should seek the simple solution.*

— Anton Pavlovich Chekhov (c. 1890)

*I love deadlines. I like the whooshing sound they make as they fly by.*

— Douglas Adams

## \*8 Matroids

### 8.1 Definitions

Many problems that can be correctly solved by greedy algorithms can be described in terms of an abstract combinatorial object called a *matroid*. Matroids were first described in 1935 by the mathematician Hassler Whitney as a combinatorial generalization of linear independence of vectors—‘matroid’ means ‘something sort of like a matrix’.

A matroid  $\mathcal{M}$  is a finite collection of finite sets that satisfies three axioms:

- **Non-emptiness:** The empty set is in  $\mathcal{M}$ . (Thus,  $\mathcal{M}$  is not itself empty.)
- **Heredity:** If a set  $X$  is an element of  $\mathcal{M}$ , then any subset of  $X$  is also in  $\mathcal{M}$ .
- **Exchange:** If  $X$  and  $Y$  are two sets in  $\mathcal{M}$  and  $|X| > |Y|$ , then there is an element  $x \in X \setminus Y$  such that  $Y \cup \{x\}$  is in  $\mathcal{M}$ .

The sets in  $\mathcal{M}$  are typically called *independent sets*; for example, we would say that any subset of an independent set is independent. The union of all sets in  $\mathcal{M}$  is called the *ground set*. An independent set is called a *basis* if it is not a proper subset of another independent set. The exchange property implies that every basis of a matroid has the same cardinality. The *rank* of a subset  $X$  of the ground set is the size of the largest independent subset of  $X$ . A subset of the ground set that is not in  $\mathcal{M}$  is called *dependent* (surprise, surprise). Finally, a dependent set is called a *circuit* if every proper subset is independent.

Most of this terminology is justified by Whitney’s original example:

**Linear matroid:** Let  $A$  be any  $n \times m$  matrix. A subset  $I \subseteq \{1, 2, \dots, n\}$  is independent if and only if the corresponding subset of columns of  $A$  is linearly independent.

The heredity property follows directly from the definition of linear independence; the exchange property is implied by an easy dimensionality argument. A basis in any linear matroid is also a basis (in the linear-algebra sense) of the vector space spanned by the columns of  $A$ . Similarly, the rank of a set of indices is precisely the rank (in the linear-algebra sense) of the corresponding set of column vectors.

Here are several other examples of matroids; some of these we will see again later. I will leave the proofs that these are actually matroids as exercises for the reader.

- **Uniform matroid  $U_{k,n}$ :** A subset  $X \subseteq \{1, 2, \dots, n\}$  is independent if and only if  $|X| \leq k$ . Any subset of  $\{1, 2, \dots, n\}$  of size  $k$  is a basis; any subset of size  $k+1$  is a circuit.
- **Graphic/cycle matroid  $\mathcal{M}(G)$ :** Let  $G = (V, E)$  be an arbitrary undirected graph. A subset of  $E$  is independent if it defines an acyclic subgraph of  $G$ . A basis in the graphic matroid is a spanning tree of  $G$ ; a circuit in this matroid is a cycle in  $G$ .

- **Cographic/cocycle matroid  $\mathcal{M}^*(G)$ :** Let  $G = (V, E)$  be an arbitrary undirected graph. A subset  $I \subseteq E$  is independent if the complementary subgraph  $(V, E \setminus I)$  of  $G$  is *connected*. A basis in this matroid is the complement of a spanning tree; a circuit in this matroid is a *cocycle*—a minimal set of edges that disconnects the graph.
- **Matching matroid:** Let  $G = (V, E)$  be an arbitrary undirected graph. A subset  $I \subseteq V$  is independent if there is a matching in  $G$  that covers  $I$ .
- **Disjoint path matroid:** Let  $G = (V, E)$  be an arbitrary *directed* graph, and let  $s$  be a fixed vertex of  $G$ . A subset  $I \subseteq V$  is independent if and only if there are edge-disjoint paths from  $s$  to each vertex in  $I$ .

Now suppose each element of the ground set of a matroid  $\mathcal{M}$  is given an arbitrary non-negative weight. The **matroid optimization problem** is to compute a basis with maximum total weight. For example, if  $\mathcal{M}$  is the cycle matroid for a graph  $G$ , the matroid optimization problem asks us to find the maximum spanning tree of  $G$ . Similarly, if  $\mathcal{M}$  is the cocycle matroid for  $G$ , the matroid optimization problem seeks (the complement of) the *minimum* spanning tree.

The following natural greedy strategy computes a basis for any weighted matroid:

```

GREEDYBASIS( $\mathcal{M}, w$ ):
   $X[1..n] \leftarrow \bigcup \mathcal{M}$  (the ground set)
  sort  $X$  in decreasing order of weight  $w$ 
   $G \leftarrow \emptyset$ 
  for  $i \leftarrow 1$  to  $n$ 
    if  $G \cup \{X[i]\} \in \mathcal{M}$ 
      add  $X[i]$  to  $G$ 
  return  $G$ 

```

Suppose we can test in  $F(n)$  whether a given subset of the ground set is independent. Then this algorithm runs in  $O(n \log n + n \cdot F(n))$  time.

**Theorem 1.** *For any matroid  $\mathcal{M}$  and any weight function  $w$ ,  $\text{GREEDYBASIS}(\mathcal{M}, w)$  returns a maximum-weight basis of  $\mathcal{M}$ .*

**Proof:** Let  $G = \{g_1, g_2, \dots, g_k\}$  be the independent set returned by  $\text{GREEDYBASIS}(\mathcal{M}, w)$ . If any other element could be added to  $G$  to get a larger independent set, the greedy algorithm would have added it. Thus,  $G$  is a basis.

For purposes of deriving a contradiction, suppose there is an independent set  $H = \{h_1, h_2, \dots, h_\ell\}$  such that

$$\sum_{i=1}^k w(g_i) < \sum_{j=1}^\ell w(h_j).$$

Without loss of generality, we assume that  $H$  is a basis. The exchange property now implies that  $k = \ell$ .

Now suppose the elements of  $G$  and  $H$  are indexed in order of decreasing weight. Let  $i$  be the smallest index such that  $w(g_i) < w(h_i)$ , and consider the independent sets

$$G_{i-1} = \{g_1, g_2, \dots, g_{i-1}\} \quad \text{and} \quad H_i = \{h_1, h_2, \dots, h_{i-1}, h_i\}.$$

By the exchange property, there is some element  $h_j \in H_i$  such that  $G_{i-1} \cup \{h_j\}$  is an independent set. We have  $w(h_j) \geq w(h_i) > w(g_i)$ . Thus, the greedy algorithm considers *and rejects* the heavier element  $h_j$  before it considers the lighter element  $g_i$ . But this is impossible—the greedy algorithm accepts elements in decreasing order of weight.  $\square$

We now immediately have a correct greedy optimization algorithm for *any* matroid. Returning to our examples:

- Linear matroid: Given a matrix  $A$ , compute a subset of vectors of maximum total weight that span the column space of  $A$ .
- Uniform matroid: Given a set of weighted objects, compute its  $k$  largest elements.
- Cycle matroid: Given a graph with weighted edges, compute its maximum spanning tree. In this setting, the greedy algorithm is better known as *Kruskal's algorithm*.
- Cocycle matroid: Given a graph with weighted edges, compute its minimum spanning tree.
- Matching matroid: Given a graph, determine whether it has a perfect matching.
- Disjoint path matroid: Given a directed graph with a special vertex  $s$ , find the largest set of edge-disjoint paths from  $s$  to other vertices.

The exchange condition for matroids turns out to be crucial for the success of this algorithm. A *subset system* is a finite collection  $\mathcal{S}$  of finite sets that satisfies the heredity condition—if  $X \in \mathcal{S}$  and  $Y \subseteq X$ , then  $Y \in \mathcal{S}$ —but not necessarily the exchange condition.

**Theorem 2.** *For any subset system  $\mathcal{S}$  that is **not** a matroid, there is a weight function  $w$  such that  $\text{GREEDYBASIS}(\mathcal{S}, w)$  does **not** return a maximum-weight set in  $\mathcal{S}$ .*

**Proof:** Let  $X$  and  $Y$  be two sets in  $\mathcal{S}$  that violate the exchange property— $|X| > |Y|$ , but for any element  $x \in X \setminus Y$ , the set  $Y \cup \{x\}$  is not in  $\mathcal{S}$ . Let  $m = |Y|$ . We define a weight function as follows:

- Every element of  $Y$  has weight  $m + 2$ .
- Every element of  $X \setminus Y$  has weight  $m + 1$ .
- Every other element of the ground set has weight zero.

With these weights, the greedy algorithm will consider and accept every element of  $Y$ , then consider and reject every element of  $X$ , and finally consider all the other elements. The algorithm returns a set with total weight  $m(m + 2) = m^2 + 2m$ . But the total weight of  $X$  is at least  $(m + 1)^2 = m^2 + 2m + 1$ . Thus, the output of the greedy algorithm is not the maximum-weight set in  $\mathcal{S}$ .  $\square$

Recall the Applied Chaos scheduling problem considered in the previous lecture note. There is a natural subset system associated with this problem: A set of classes is independent if and only if not two classes overlap. (This is just the graph-theory notion of ‘independent set’!) This subset system is *not* a matroid, because there can be maximal independent sets of different sizes, which violates the exchange property. If we consider a *weighted* version of the class scheduling problem, say where each class is worth a different number of hours, Theorem 2 implies that the greedy algorithm will *not* always find the optimal schedule. (In fact, there’s an easy counterexample with only two classes!) However, Theorem 2 does *not* contradict the correctness of the greedy algorithm for the original *unweighted* problem, however; that problem uses a particularly lucky choice of weights (all equal).

## 8.2 Scheduling with Deadlines

Suppose you have  $n$  tasks to complete in  $n$  days; each task requires your attention for a full day. Each task comes with a *deadline*, the last day by which the job should be completed, and a *penalty* that you must pay if you do not complete each task by its assigned deadline. What order should you perform your tasks in to minimize the total penalty you must pay?

More formally, you are given an array  $D[1..n]$  of deadlines and an array  $P[1..n]$  of penalties. Each deadline  $D[i]$  is an integer between 1 and  $n$ , and each penalty  $P[i]$  is a non-negative real number. A *schedule* is a permutation of the integers  $\{1, 2, \dots, n\}$ . The scheduling problem asks you to find a schedule  $\pi$  that minimizes the following cost:

$$\text{cost}(\pi) := \sum_{i=1}^n P[i] \cdot [\pi(i) > D[i]].$$

This doesn't look anything like a matroid optimization problem. For one thing, matroid optimization problems ask us to find an optimal set; this problem asks us to find an optimal *permutation*. Surprisingly, however, this scheduling problem is actually a matroid optimization in disguise! For any schedule  $\pi$ , call tasks  $i$  such that  $\pi(i) > D[i]$  *late*, and all other tasks *on time*. The following trivial observation is the key to revealing the underlying matroid structure.

The cost of a schedule is determined by the subset of tasks that are on time.

Call a subset  $X$  of the tasks *realistic* if there is a schedule  $\pi$  in which every task in  $X$  is on time. We can precisely characterize the realistic subsets as follows. Let  $X(t)$  denote the subset of tasks in  $X$  whose deadline is on or before  $t$ :

$$X(t) := \{i \in X \mid D[i] \leq t\}.$$

In particular,  $X(0) = \emptyset$  and  $X(n) = X$ .

**Lemma 3.** *Let  $X \subseteq \{1, 2, \dots, n\}$  be an arbitrary subset of the  $n$  tasks.  $X$  is realistic if and only if  $|X(t)| \leq t$  for every integer  $t$ .*

**Proof:** Let  $\pi$  be a schedule in which every task in  $X$  is on time. Let  $i_t$  be the  $t$ th task in  $X$  to be completed. On the one hand, we have  $\pi(i_t) \geq t$ , since otherwise, we could not have completed  $t - 1$  other jobs in  $X$  before  $i_t$ . On the other hand,  $\pi(i_t) \leq D[i]$ , because  $i_t$  is on time. We conclude that  $D[i_t] \geq t$ , which immediately implies that  $|X(t)| \leq t$ .

Now suppose  $|X(t)| \leq t$  for every integer  $t$ . If we perform the tasks in  $X$  in increasing order of deadline, then we complete all tasks in  $X$  with deadlines  $t$  or less by day  $t$ . In particular, for any  $i \in X$ , we perform task  $i$  on or before its deadline  $D[i]$ . Thus,  $X$  is realistic.  $\square$

We can define a *canonical schedule* for any set  $X$  as follows: execute the tasks in  $X$  in increasing deadline order, and then execute the remaining tasks in any order. The previous proof implies that a set  $X$  is realistic if and only if every task in  $X$  is on time in the canonical schedule for  $X$ . Thus, our scheduling problem can be rephrased as follows:

Find a realistic subset  $X$  such that  $\sum_{i \in X} P[i]$  is maximized.

So we're looking for optimal subsets after all.

**Lemma 4.** *The collection of realistic sets of jobs forms a matroid.*

**Proof:** The empty set is vacuously realistic, and any subset of a realistic set is clearly realistic. Thus, to prove the lemma, it suffices to show that the exchange property holds. Let  $X$  and  $Y$  be realistic sets of jobs with  $|X| > |Y|$ .

Let  $t^*$  be the largest integer such that  $|X(t^*)| \leq |Y(t^*)|$ . This integer must exist, because  $|X(0)| = 0 \leq 0 = |Y(0)|$  and  $|X(n)| = |X| > |Y| = |Y(n)|$ . By definition of  $t^*$ , there are more tasks with deadline  $t^* + 1$  in  $X$  than in  $Y$ . Thus, we can choose a task  $j$  in  $X \setminus Y$  with deadline  $t^* + 1$ ; let  $Z = Y \cup \{j\}$ .

Let  $t$  be an arbitrary integer. If  $t \leq t^*$ , then  $|Z(t)| = |Y(t)| \leq t$ , because  $Y$  is realistic. On the other hand, if  $t > t^*$ , then  $|Z(t)| = |Y(t)| + 1 \leq |X(t)| < t$  by definition of  $t^*$  and because  $X$  is realistic. The previous lemma now implies that  $Z$  is realistic. This completes the proof of the exchange property.  $\square$

This lemma implies that our scheduling problem is a matroid optimization problem, so the greedy algorithm finds the optimal schedule.

```

GREEDYSCHEDULE( $D[1..n], P[1..n]$ ):
    Sort  $P$  in increasing order, and permute  $D$  to match
     $j \leftarrow 0$ 
    for  $i \leftarrow 1$  to  $n$ 
         $X[j+1] \leftarrow i$ 
        if  $X[1..j+1]$  is realistic
             $j \leftarrow j+1$ 
    return the canonical schedule for  $X[1..j]$ 

```

To turn this outline into a real algorithm, we need a procedure to test whether a given subset of jobs is realistic. Lemma 9 immediately suggests the following strategy to answer this question in  $O(n)$  time.

```

REALISTIC?( $X[1..m], D[1..n]$ ):
    « $X$  is sorted by increasing deadline:  $i \leq j \implies D[X[i]] \leq D[X[j]]$ »
     $N \leftarrow 0$ 
     $j \leftarrow 0$ 
    for  $t \leftarrow 1$  to  $n$ 
        if  $D[X[j]] = t$ 
             $N \leftarrow N + 1$ ;  $j \leftarrow j + 1$ 
    «Now  $N = |X(t)|$ »
    if  $N > t$ 
        return FALSE
    return TRUE

```

If we use this subroutine, GREEDYSCHEDULE runs in  $O(n^2)$  time. By using some appropriate data structures, the running time can be reduced to  $O(n \log n)$ ; details are left as an exercise for the reader.

## Exercises

1. Prove that for any graph  $G$ , the ‘graphic matroid’  $\mathcal{M}(G)$  is in fact a matroid. (This problem is really asking you to prove that Kruskal’s algorithm is correct!)
2. Prove that for any graph  $G$ , the ‘cographic matroid’  $\mathcal{M}^*(G)$  is in fact a matroid.
3. Prove that for any graph  $G$ , the ‘matching matroid’ of  $G$  is in fact a matroid. [*Hint: What is the symmetric difference of two matchings?*]

4. Prove that for any directed graph  $G$  and any vertex  $s$  of  $G$ , the resulting ‘disjoint path matroid’ of  $G$  is in fact a matroid. [*Hint: This question is **much** easier if you’re already familiar with maximum flows.*]
5. Let  $G$  be an undirected graph. A set of cycles  $\{c_1, c_2, \dots, c_k\}$  in  $G$  is called *redundant* if every edge in  $G$  appears in an even number of  $c_i$ ’s. A set of cycles is *independent* if it contains no redundant subset. A maximal independent set of cycles is called a *cycle basis* for  $G$ .
  - (a) Let  $C$  be any cycle basis for  $G$ . Prove that for any cycle  $\gamma$  in  $G$ , there is a subset  $A \subseteq C$  such that  $A \cap \{\gamma\}$  is redundant. In other words,  $\gamma$  is the ‘exclusive or’ of the cycles in  $A$ .
  - (b) Prove that the set of independent cycle sets form a matroid.

\*(c) Now suppose each edge of  $G$  has a weight. Define the weight of a cycle to be the total weight of its edges, and the weight of a set of cycles to be the total weight of all cycles in the set. (Thus, each edge is counted once for every cycle in which it appears.) Describe and analyze an efficient algorithm to compute the minimum-weight cycle basis in  $G$ .
6. Describe a modification of GREEDYSCHEDULE that runs in  $O(n \log n)$  time. [*Hint: Store  $X$  in an appropriate data structure that supports the operations “Is  $X \cup \{i\}$  realistic?” and “Add  $i$  to  $X$ ” in  $O(\log n)$  time each.*]

*The first nuts and bolts appeared in the middle 1400's. The bolts were just screws with straight sides and a blunt end. The nuts were hand-made, and very crude. When a match was found between a nut and a bolt, they were kept together until they were finally assembled.*

*In the Industrial Revolution, it soon became obvious that threaded fasteners made it easier to assemble products, and they also meant more reliable products. But the next big step came in 1801, with Eli Whitney, the inventor of the cotton gin. The lathe had been recently improved. Batches of bolts could now be cut on different lathes, and they would all fit the same nut.*

*Whitney set up a demonstration for President Adams, and Vice-President Jefferson. He had piles of musket parts on a table. There were 10 similar parts in each pile. He went from pile to pile, picking up a part at random. Using these completely random parts, he quickly put together a working musket.*

— Karl S. Kruszelnicki ('Dr. Karl'), *Karl Trek*, December 1997

*Dr [John von] Neumann in his Theory of Games and Economic Behavior introduces the cut-up method of random action into game and military strategy: Assume that the worst has happened and act accordingly. If your strategy is at some point determined... by random factor your opponent will gain no advantage from knowing your strategy since he cannot predict the move. The cut-up method could be used to advantage in processing scientific data. How many discoveries have been made by accident? We cannot produce accidents to order.*

— William S. Burroughs, "The Cut-Up Method of Brion Gysin"  
in *The Third Mind* by William S. Burroughs and Brion Gysin (1978)

## 9 Randomized Algorithms

### 9.1 Nuts and Bolts

Suppose we are given  $n$  nuts and  $n$  bolts of different sizes. Each nut matches exactly one bolt and vice versa. The nuts and bolts are all almost exactly the same size, so we can't tell if one bolt is bigger than the other, or if one nut is bigger than the other. If we try to match a nut with a bolt, however, the nut will be either too big, too small, or just right for the bolt.

Our task is to match each nut to its corresponding bolt. But before we do this, let's try to solve some simpler problems, just to get a feel for what we can and can't do.

Suppose we want to find the nut that matches a particular bolt. The obvious algorithm — test every nut until we find a match — requires exactly  $n - 1$  tests in the worst case. We might have to check every bolt except one; if we get down to the last bolt without finding a match, we know that the last nut is the one we're looking for.<sup>1</sup>

Intuitively, in the 'average' case, this algorithm will look at approximately  $n/2$  nuts. But what exactly does 'average case' mean?

### 9.2 Deterministic vs. Randomized Algorithms

Normally, when we talk about the running time of an algorithm, we mean the *worst-case* running time. This is the maximum, over all problems of a certain size, of the running time of that algorithm on that input:

$$T_{\text{worst-case}}(n) = \max_{|X|=n} T(X).$$

On extremely rare occasions, we will also be interested in the *best-case* running time:

$$T_{\text{best-case}}(n) = \min_{|X|=n} T(X).$$

---

<sup>1</sup>"Whenever you lose something, it's always in the last place you look. So why not just look there first?"

The *average-case* running time is best defined by the *expected value*, over all inputs  $X$  of a certain size, of the algorithm's running time for  $X$ :<sup>2</sup>

$$T_{\text{average-case}}(n) = \underset{|X|=n}{\mathbb{E}}[T(X)] = \sum_{|X|=n} T(x) \cdot \Pr[X].$$

The problem with this definition is that we rarely, if ever, know what the probability of getting any particular input  $X$  is. We could compute average-case running times by assuming a particular probability distribution—for example, every possible input is equally likely—but this assumption doesn't describe reality very well. Most real-life data is decidedly non-random (or at least random in some unpredictable way).

Instead of considering this rather questionable notion of average case running time, we will make a distinction between two kinds of algorithms: *deterministic* and *randomized*. A deterministic algorithm is one that always behaves the same way given the same input; the input completely *determines* the sequence of computations performed by the algorithm. Randomized algorithms, on the other hand, base their behavior not only on the input but also on several *random* choices. The same randomized algorithm, given the same input multiple times, may perform different computations in each invocation.

This means, among other things, that the running time of a randomized algorithm on a given input is no longer fixed, but is itself a random variable. When we analyze randomized algorithms, we are typically interested in the *worst-case expected* running time. That is, we look at the average running time for each input, and then choose the maximum over all inputs of a certain size:

$$T_{\text{worst-case expected}}(n) = \max_{|X|=n} \mathbb{E}[T(X)].$$

It's important to note here that we are making *no* assumptions about the probability distribution of possible inputs. All the randomness is inside the algorithm, where we can control it!

### 9.3 Back to Nuts and Bolts

Let's go back to the problem of finding the nut that matches a given bolt. Suppose we use the same algorithm as before, but at each step we choose a nut *uniformly at random* from the untested nuts. 'Uniformly' is a technical term meaning that each nut has exactly the same probability of being chosen.<sup>3</sup> So if there are  $k$  nuts left to test, each one will be chosen with probability  $1/k$ . Now what's the expected number of comparisons we have to perform? Intuitively, it should be about  $n/2$ , but let's formalize our intuition.

Let  $T(n)$  denote the number of comparisons our algorithm uses to find a match for a single bolt out of  $n$  nuts.<sup>4</sup> We still have some simple base cases  $T(1) = 0$  and  $T(2) = 1$ , but when  $n > 2$ ,  $T(n)$  is a random variable.  $T(n)$  is always between 1 and  $n - 1$ ; its actual value depends on our algorithm's random choices. We are interested in the *expected value* or *expectation* of  $T(n)$ , which is defined as follows:

$$\mathbb{E}[T(n)] = \sum_{k=1}^{n-1} k \cdot \Pr[T(n) = k]$$

---

<sup>2</sup>The notation  $\mathbb{E}[\cdot]$  for expectation has nothing to do with the shift operator  $\mathbf{E}$  used in the annihilator method for solving recurrences!

<sup>3</sup>This is what most people think 'random' means, but they're wrong.

<sup>4</sup>Note that for this algorithm, the input is completely specified by the number  $n$ . Since we're choosing the nuts to test at random, even the order in which the nuts and bolts are presented doesn't matter. That's why I'm using the simpler notation  $T(n)$  instead of  $T(X)$ .

If the target nut is the  $k$ th nut tested, our algorithm performs  $\min\{k, n-1\}$  comparisons. In particular, if the target nut is the last nut chosen, we don't actually test it. Because we choose the next nut to test uniformly at random, the target nut is equally likely—with probability exactly  $1/n$ —to be the first, second, third, or  $k$ th bolt tested, for any  $k$ . Thus:

$$\Pr[T(n) = k] = \begin{cases} 1/n & \text{if } k < n-1, \\ 2/n & \text{if } k = n-1. \end{cases}$$

Plugging this into the definition of expectation gives us our answer.

$$\begin{aligned} E[T(n)] &= \sum_{k=1}^{n-2} \frac{k}{n} + \frac{2(n-1)}{n} \\ &= \sum_{k=1}^{n-1} \frac{k}{n} + \frac{n-1}{n} \\ &= \frac{n(n-1)}{2n} + 1 - \frac{1}{n} \\ &= \frac{n+1}{2} - \frac{1}{n} \end{aligned}$$

We can get exactly the same answer by thinking of this algorithm recursively. We always have to perform at least one test. With probability  $1/n$ , we successfully find the matching nut and halt. With the remaining probability  $1 - 1/n$ , we recursively solve the same problem but with one fewer nut. We get the following recurrence for the expected number of tests:

$$T(1) = 0, \quad E[T(n)] = 1 + \frac{n-1}{n} E[T(n-1)]$$

To get the solution, we define a new function  $t(n) = n E[T(n)]$  and rewrite:

$$t(1) = 0, \quad t(n) = n + t(n-1)$$

This recurrence translates into a simple summation, which we can easily solve.

$$\begin{aligned} t(n) &= \sum_{k=2}^n k = \frac{n(n+1)}{2} - 1 \\ \implies E[T(n)] &= \frac{t(n)}{n} = \frac{n+1}{2} - \frac{1}{n} \end{aligned}$$

## 9.4 Finding All Matches

Not let's go back to the problem introduced at the beginning of the lecture: finding the matching nut for every bolt. The simplest algorithm simply compares every nut with every bolt, for a total of  $n^2$  comparisons. The next thing we might try is repeatedly finding an arbitrary matched pair, using our very first nuts and bolts algorithm. This requires

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

comparisons in the worst case. So we save roughly a factor of two over the really stupid algorithm. Not very exciting.

Here's another possibility. Choose a *pivot* bolt, and test it against every nut. Then test the matching pivot nut against every other bolt. After these  $2n - 1$  tests, we have one matched pair, and the remaining nuts and bolts are partitioned into two subsets: those smaller than the pivot pair and those larger than the pivot pair. Finally, recursively match up the two subsets. The worst-case number of tests made by this algorithm is given by the recurrence

$$\begin{aligned} T(n) &= 2n - 1 + \max_{1 \leq k \leq n} \{T(k - 1) + T(n - k)\} \\ &= 2n - 1 + T(n - 1) \end{aligned}$$

Along with the trivial base case  $T(0) = 0$ , this recurrence solves to

$$T(n) = \sum_{i=1}^n (2n - 1) = n^2.$$

In the worst case, this algorithm tests *every* nut-bolt pair! We could have been a little more clever—for example, if the pivot bolt is the smallest bolt, we only need  $n - 1$  tests to partition everything, not  $2n - 1$ —but cleverness doesn't actually help that much. We still end up with about  $n^2/2$  tests in the worst case.

However, since this recursive algorithm looks almost exactly like quicksort, and everybody 'knows' that the 'average-case' running time of quicksort is  $\Theta(n \log n)$ , it seems reasonable to guess that the average number of nut-bolt comparisons is also  $\Theta(n \log n)$ . As we shall see shortly, if the pivot bolt is always chosen *uniformly at random*, this intuition is exactly right.

## 9.5 Reductions to and from Sorting

The second algorithm for matching up the nuts and bolts looks exactly like quicksort. The algorithm not only matches up the nuts and bolts, but also sorts them by size.

In fact, the problems of sorting and matching nuts and bolts are equivalent, in the following sense. If the bolts were sorted, we could match the nuts and bolts in  $O(n \log n)$  time by performing a binary search with each nut. Thus, if we had an algorithm to sort the bolts in  $O(n \log n)$  time, we would immediately have an algorithm to match the nuts and bolts, starting from scratch, in  $O(n \log n)$  time. This process of *assuming* a solution to one problem and using it to solve another is called *reduction*—we can *reduce* the matching problem to the sorting problem in  $O(n \log n)$  time.

There is a reduction in the other direction, too. If the nuts and bolts were matched, we could sort them in  $O(n \log n)$  time using, for example, merge sort. Thus, if we have an  $O(n \log n)$  time algorithm for either sorting or matching nuts and bolts, we automatically have an  $O(n \log n)$  time algorithm for the other problem.

Unfortunately, since we aren't allowed to directly compare two bolts or two nuts, we can't use heapsort or mergesort to sort the nuts and bolts in  $O(n \log n)$  worst case time. In fact, the problem of sorting nuts and bolts *deterministically* in  $O(n \log n)$  time was only 'solved' in 1995<sup>5</sup>, but both the algorithms and their analysis are incredibly technical and the constant hidden in the  $O(\cdot)$  notation is quite large.

Reductions will come up again later in the course when we start talking about lower bounds and NP-completeness.

---

<sup>5</sup>János Komlós, Yuan Ma, and Endre Szemerédi, Sorting nuts and bolts in  $O(n \log n)$  time, *SIAM J. Discrete Math* 11(3):347–372, 1998. See also Phillip G. Bradford, Matching nuts and bolts optimally, Technical Report MPI-I-95-1-025, Max-Planck-Institut für Informatik, September 1995. Bradford's algorithm is *slightly* simpler.

## 9.6 Recursive Analysis

Intuitively, we can argue that our quicksort-like algorithm will usually choose a bolt of approximately median size, and so the average numbers of tests should be  $O(n \log n)$ . We can now finally formalize this intuition. To simplify the notation slightly, I'll write  $\bar{T}(n)$  in place of  $E[T(n)]$  everywhere.

Our randomized matching/sorting algorithm chooses its pivot bolt *uniformly at random* from the set of unmatched bolts. Since the pivot bolt is equally likely to be the smallest, second smallest, or  $k$ th smallest for any  $k$ , the expected number of tests performed by our algorithm is given by the following recurrence:

$$\begin{aligned}\bar{T}(n) &= 2n - 1 + E_k [\bar{T}(k-1) + \bar{T}(n-k)] \\ &= 2n - 1 + \frac{1}{n} \sum_{k=1}^n (\bar{T}(k-1) + \bar{T}(n-k))\end{aligned}$$

The base case is  $T(0) = 0$ . (We can save a few tests by setting  $T(1) = 0$  instead of 1, but the analysis will be easier if we're a little stupid.)

Yuck. At this point, we could simply *guess* the solution, based on the incessant rumors that quicksort runs in  $O(n \log n)$  time in the average case, and prove our guess correct by induction. (See Section 9.8 below for details.)

However, if we're only interested in asymptotic bounds, we can afford to be a little conservative. What we'd *really* like is for the pivot bolt to be the median bolt, so that half the bolts are bigger and half the bolts are smaller. This isn't very likely, but there is a good chance that the pivot bolt is close to the median bolt. Let's say that a pivot bolt is *good* if it's in the middle half of the final sorted set of bolts, that is, bigger than at least  $n/4$  bolts and smaller than at least  $n/4$  bolts. If the pivot bolt is good, then the *worst* split we can have is into one set of  $3n/4$  pairs and one set of  $n/4$  pairs. If the pivot bolt is bad, then our algorithm is still better than starting over from scratch. Finally, a randomly chosen pivot bolt is good with probability  $1/2$ .

These simple observations give us the following simple recursive *upper bound* for the expected running time of our algorithm:

$$\bar{T}(n) \leq 2n - 1 + \frac{1}{2} \left( \bar{T}\left(\frac{3n}{4}\right) + \bar{T}\left(\frac{n}{4}\right) \right) + \frac{1}{2} \cdot \bar{T}(n)$$

A little algebra simplifies this even further:

$$\bar{T}(n) \leq 4n - 2 + \bar{T}\left(\frac{3n}{4}\right) + \bar{T}\left(\frac{n}{4}\right)$$

We can solve this recurrence using the recursion tree method, giving us the unsurprising upper bound  $\bar{T}(n) = O(n \log n)$ . A similar argument gives us the matching lower bound  $\bar{T}(n) = \Omega(n \log n)$ .

Unfortunately, while this argument is convincing, it is *not* a formal proof, because it relies on the unproven assumption that  $\bar{T}(n)$  is a *convex* function, which means that  $\bar{T}(n+1) + \bar{T}(n-1) \geq 2\bar{T}(n)$  for all  $n$ .  $\bar{T}(n)$  is actually convex, but we never proved it. Convexity follows from the closed-form solution of the recurrence, but using that fact would be circular logic. Sadly, formally proving convexity seems to be almost as hard as solving the recurrence. If we want a *proof* of the expected cost of our algorithm, we need another way to proceed.

## 9.7 Iterative Analysis

By making a simple change to our algorithm, which has no effect on the number of tests, we can analyze it much more directly and exactly, without solving a recurrence or relying on hand-wavy intuition.

The recursive subproblems solved by quicksort can be laid out in a binary tree, where each node corresponds to a subset of the nuts and bolts. In the usual recursive formulation, the algorithm partitions the nuts and bolts at the root, then the left child of the root, then the leftmost grandchild, and so forth, recursively sorting everything on the left before starting on the right subproblem.

But we don't have to solve the subproblems in this order. In fact, we can visit the nodes in the recursion tree in any order we like, as long as the root is visited first, and any other node is visited after its parent. Thus, we can recast quicksort in the following iterative form. Choose a pivot bolt, find its match, and partition the remaining nuts and bolts into two subsets. Then pick a second pivot bolt and partition whichever of the two subsets contains it. At this point, we have two matched pairs and three subsets of nuts and bolts. Continue choosing new pivot bolts and partitioning subsets, each time finding one match and increasing the number of subsets by one, until every bolt has been chosen as the pivot. At the end, every bolt has been matched, and the nuts and bolts are sorted.

Suppose we always choose the next pivot bolt *uniformly at random* from the bolts that haven't been pivots yet. Then no matter which subset contains this bolt, the pivot bolt is equally likely to be any bolt *in that subset*. That implies (by induction) that our randomized iterative algorithm performs *exactly* the same set of tests as our randomized recursive algorithm, but possibly in a different order.

Now let  $B_i$  denote the  $i$ th smallest bolt, and  $N_j$  denote the  $j$ th smallest nut. For each  $i$  and  $j$ , define an indicator variable  $X_{ij}$  that equals 1 if our algorithm compares  $B_i$  with  $N_j$  and zero otherwise. Then the total number of nut/bolt comparisons is exactly

$$T(n) = \sum_{i=1}^n \sum_{j=1}^n X_{ij}.$$

We are interested in the expected value of this double summation:

$$\mathbb{E}[T(n)] = \mathbb{E}\left[\sum_{i=1}^n \sum_{j=1}^n X_{ij}\right] = \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}[X_{ij}].$$

This equation uses a crucial property of random variables called *linearity of expectation*: for any random variables  $X$  and  $Y$ , the sum of their expectations is equal to the expectation of their sum:  $\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$ .

To analyze our algorithm, we only need to compute the expected value of each  $X_{ij}$ . By definition of expectation,

$$\mathbb{E}[X_{ij}] = 0 \cdot \Pr[X_{ij} = 0] + 1 \cdot \Pr[X_{ij} = 1] = \Pr[X_{ij} = 1],$$

so we just need to calculate  $\Pr[X_{ij} = 1]$  for all  $i$  and  $j$ .

First let's assume that  $i < j$ . The only comparisons our algorithm performs are between some pivot bolt (or its partner) and a nut (or bolt) in the same subset. The only event that can prevent a comparison between  $B_i$  and  $N_j$  is choosing some intermediate pivot bolt  $B_k$ , with  $i < k < j$ , before  $B_i$  or  $B_j$ . In other words:

**Our algorithm compares  $B_i$  and  $N_j$  if and only if the first pivot chosen from the set  $\{B_i, B_{i+1}, \dots, B_j\}$  is either  $B_i$  or  $B_j$ .**

Since the set  $\{B_i, B_{i+1}, \dots, B_j\}$  contains  $j - i + 1$  bolts, each of which is equally likely to be chosen first, we immediately have

$$\mathbb{E}[X_{ij}] = \frac{2}{j - i + 1} \quad \text{for all } i < j.$$

Symmetric arguments give us  $E[X_{ij}] = \frac{2}{i-j+1}$  for all  $i > j$ . Since our algorithm is a little stupid, every bolt is compared with its partner, so  $X_{ii} = 1$  for all  $i$ . (In fact, if a pivot bolt is the only bolt in its subset, we don't need to compare it against its partner, but this improvement complicates the analysis.)

Putting everything together, we get the following summation.

$$\begin{aligned} E[T(n)] &= \sum_{i=1}^n \sum_{j=1}^n E[X_{ij}] \\ &= \sum_{i=1}^n E[X_{ii}] + 2 \sum_{i=1}^n \sum_{j=i+1}^n E[X_{ij}] \\ &= \boxed{n + 4 \sum_{i=1}^n \sum_{j=i+1}^n \frac{1}{j-i+1}} \end{aligned}$$

This is quite a bit simpler than the recurrence we got before. With just a few more lines of algebra, we can turn it into an exact, closed-form expression for the expected number of comparisons.

$$\begin{aligned} E[T(n)] &= n + 4 \sum_{i=1}^n \sum_{k=2}^{n-i+1} \frac{1}{k} && [\text{substitute } k = j - i + 1] \\ &= n + 4 \sum_{k=2}^n \sum_{i=1}^{n-k+1} \frac{1}{k} && [\text{reorder summations}] \\ &= n + 4 \sum_{k=2}^n \frac{n-k+1}{k} \\ &= n + 4 \left( (n+1) \sum_{k=2}^n \frac{1}{k} - \sum_{k=2}^n 1 \right) \\ &= n + 4((n+1)(H_n - 1) - (n-1)) \\ &= n + 4(nH_n - 2n + H_n) \\ &= \boxed{4nH_n - 7n + 4H_n} \end{aligned}$$

Sure enough, it's  $\Theta(n \log n)$ .

## \*9.8 Masochistic Analysis

If we're feeling particularly masochistic, we can actually solve the recurrence directly, all the way to an exact closed-form solution. I'm including this only to show you it can be done; this won't be on the test.

First we simplify the recurrence slightly by combining symmetric terms.

$$\bar{T}(n) = \frac{1}{n} \sum_{k=1}^n (\bar{T}(k-1) + \bar{T}(n-k)) + 2n - 1 = \frac{2}{n} \sum_{k=0}^{n-1} \bar{T}(k) + 2n - 1$$

We then convert this 'full history' recurrence into a 'limited history' recurrence by shifting and subtracting away common terms. (I call this "Magic step #1".) To make this step slightly easier, we first multiply

both sides of the recurrence by  $n$  to get rid of the fractions.

$$\begin{aligned}
 n\bar{T}(n) &= 2 \sum_{k=0}^{n-1} \bar{T}(k) + 2n^2 - n \\
 (n-1)\bar{T}(n-1) &= 2 \sum_{k=0}^{n-2} \bar{T}(k) + 2(n-1)^2 - (n-1) \\
 &= 2 \sum_{k=0}^{n-2} \bar{T}(k) + 2n^2 - 5n + 3 \\
 n\bar{T}(n) - (n-1)\bar{T}(n-1) &= 2T(n-1) + 4n - 3 \\
 \bar{T}(n) &= \frac{n+1}{n}\bar{T}(n-1) + 4 - \frac{3}{n}
 \end{aligned}$$

To solve this limited-history recurrence, we define a new function  $t(n) = \bar{T}(n)/(n+1)$ . (I call this “Magic step #2”.) This gives us an even simpler recurrence for  $t(n)$  in terms of  $t(n-1)$ :

$$\begin{aligned}
 t(n) &= \frac{\bar{T}(n)}{n+1} \\
 &= \frac{1}{n+1} \left( (n+1) \frac{T(n-1)}{n} + 4 - \frac{3}{n} \right) \\
 &= t(n-1) + \frac{4}{n+1} - \frac{3}{n(n+1)} \\
 &= t(n-1) + \frac{7}{n+1} - \frac{3}{n}
 \end{aligned}$$

I used the technique of partial fractions (remember calculus?) to replace  $\frac{1}{n(n+1)}$  with  $\frac{1}{n} - \frac{1}{n+1}$  in the last step. The base case for this recurrence is  $t(0) = 0$ . Once again, we have a recurrence that translates directly into a summation, which we can solve with just a few lines of algebra.

$$\begin{aligned}
 t(n) &= \sum_{i=1}^n \left( \frac{7}{i+1} - \frac{3}{i} \right) \\
 &= 7 \sum_{i=1}^n \frac{1}{i+1} - 3 \sum_{i=1}^n \frac{1}{i} \\
 &= 7(H_{n+1} - 1) - 3H_n \\
 &= 4H_n - 7 + \frac{7}{n+1}
 \end{aligned}$$

The last step uses the recursive definition of the harmonic numbers:  $H_n = H_n + \frac{1}{n+1}$ . Finally, substituting  $\bar{T}(n) = (n+1)t(n)$  and simplifying gives us the exact solution to the original recurrence.

$$\bar{T}(n) = 4(n+1)H_n - 7(n+1) + 7 = \boxed{4nH_n - 7n + 4H_n}$$

Surprise, surprise, we get exactly the same solution!

## Exercises

### Probability

Several of these problems refer to decks of playing cards. A standard (Anglo-American) deck of 52 playing cards contains 13 cards in each of four suits: spades ( $\spadesuit$ ), hearts ( $\heartsuit$ ), diamonds ( $\diamondsuit$ ), and clubs ( $\clubsuit$ ). Within each suit, the 13 cards have distinct *ranks*: 2, 3, 4, 5, 6, 7, 8, 9, 10, jack ( $J$ ), queen ( $Q$ ), king ( $K$ ), and ace ( $A$ ). For purposes of these problems, the ranks are ordered  $A < 2 < 3 < \dots < 9 < 10 < J < Q < K$ ; thus, for example, the jack of spades has higher rank than the eight of diamonds.

1. Clock Solitaire is played with a standard deck of playing cards. To set up the game, deal the cards face down into 13 piles of four cards each, one in each of the ‘hour’ positions of a clock and one in the center. Each pile corresponds to a particular rank—A through Q in clockwise order for the hour positions, and K for the center. To start the game, turn over a card in the center pile. Then repeatedly turn over a card in the pile corresponding to the value of the previous card. The game ends when you try to turn over a card from a pile whose four cards are already face up. (This is always the center pile—why?) You win if and only if every card is face up when the game ends.

What is the *exact* probability that you win a game of Clock Solitaire, assuming that the cards are permuted uniformly at random before they are dealt into their piles?

2. Professor Jay is about to perform a public demonstration with two decks of cards, one with red backs (‘the red deck’) and one with blue backs (‘the blue deck’). Both decks lie face-down on a table in front of Professor Jay, *shuffled* so that every permutation of each deck is equally likely.

To begin the demonstration, Professor Jay turns over the top card from each deck. If one of these two cards is the three of clubs ( $3\clubsuit$ ), the demonstration ends immediately. Otherwise, the good Professor repeatedly hurls the cards he just turned over into the *thick, pachydermatous outer melon layer* of a nearby watermelon, and then turns over the next card from the top of each deck. The demonstration ends the first time a  $3\clubsuit$  is turned over. Thus, if  $3\clubsuit$  is the last card in both decks, the demonstration ends with 102 cards embedded in the watermelon, that most prodigious of household fruits.

- (a) What is the *exact* expected number of cards that Professor Jay hurls into the watermelon?
- (b) For each of the statements below, give the *exact* probability that the statement is true of the *first* pair of cards Professor Jay turns over.
  - i. Both cards are threes.
  - ii. One card is a three, and the other card is a club.
  - iii. If (at least) one card is a heart, then (at least) one card is a diamond.
  - iv. The card from the red deck has higher rank than the card from the blue deck.
- (c) For each of the statements below, give the *exact* probability that the statement is true of the *last* pair of cards Professor Jay turns over.
  - i. Both cards are threes.
  - ii. One card is a three, and the other card is a club.
  - iii. If (at least) one card is a heart, then (at least) one card is a diamond.
  - iv. The card from the red deck has higher rank than the card from the blue deck.

3. Penn and Teller agree to play the following game. Penn shuffles a standard deck of playing cards so that every permutation is equally likely. Then Teller draws cards from the deck, one at a time without replacement, until he draws the three of clubs ( $3\clubsuit$ ), at which point the remaining undrawn cards instantly burst into flames.

The first time Teller draws a card from the deck, he gives it to Penn. From then on, until the game ends, whenever Teller draws a card whose value is smaller than the last card he gave to Penn, he gives the new card to Penn.<sup>6</sup> To make the rules unambiguous, they agree beforehand that  $A = 1$ ,  $J = 11$ ,  $Q = 12$ , and  $K = 13$ .

- (a) What is the expected number of cards that Teller draws?
  - (b) What is the expected *maximum* value among the cards Teller gives to Penn?
  - (c) What is the expected *minimum* value among the cards Teller gives to Penn?
  - (d) What is the expected number of cards that Teller gives to Penn? [*Hint: Let  $13 = n$ .*]
4. Suppose  $n$  lights labeled  $0, \dots, n - 1$  are placed clockwise around a circle. Initially, every light is off. Consider the following random process.

```
LIGHTTHECIRCLE( $n$ ):
 $k \leftarrow 0$ 
turn on light 0
while at least one light is off
    with probability  $1/2$ 
         $k \leftarrow (k + 1) \bmod n$ 
    else
         $k \leftarrow (k - 1) \bmod n$ 
    if light  $k$  is off, turn it on
```

- (a) Let  $p(i, n)$  be the probability that light  $i$  is the last to be turned on by  $\text{LIGHTTHECIRCLE}(n, 0)$ . For example,  $p(0, 2) = 0$  and  $p(1, 2) = 1$ . Find an exact closed-form expression for  $p(i, n)$  in terms of  $n$  and  $i$ . Prove your answer is correct.
  - (b) Give the tightest upper bound you can on the expected running time of this algorithm.
5. Consider a random walk on a path with vertices numbered  $1, 2, \dots, n$  from left to right. At each step, we flip a coin to decide which direction to walk, moving one step left or one step right with equal probability. The random walk ends when we fall off one end of the path, either by moving left from vertex 1 or by moving right from vertex  $n$ .
- (a) Prove that the probability that the walk ends by falling off the *right* end of the path is exactly  $1/(n + 1)$ .
  - (b) Prove that if we start at vertex  $k$ , the probability that we fall off the *right* end of the path is exactly  $k/(n + 1)$ .
  - (c) Prove that if we start at vertex 1, the expected number of steps before the random walk ends is exactly  $n$ .
  - (d) Suppose we start at vertex  $n/2$  instead. State and prove a tight  $\Theta$ -bound on the expected length of the random walk in this case.

---

<sup>6</sup>Specifically, he hurls it directly into the back of Penn's right hand.

## Randomized Algorithms

6. Consider the following randomized algorithm for generating biased random bits. The subroutine FAIRCOIN returns either 0 or 1 with equal probability; the random bits returned by FAIRCOIN are mutually independent.

```

ONEINTHREE:
  if FAIRCOIN = 0
    return 0
  else
    return 1 - ONEINTHREE

```

- (a) Prove that ONEINTHREE returns 1 with probability 1/3.
  - (b) What is the *exact* expected number of times that this algorithm calls FAIRCOIN?
  - (c) Now suppose you are *given* a subroutine ONEINTHREE that generates a random bit that is equal to 1 with probability 1/3. Describe a FAIRCOIN algorithm that returns either 0 or 1 with equal probability, using ONEINTHREE as your only source of randomness.
  - (d) What is the *exact* expected number of times that your FAIRCOIN algorithm calls ONEINTHREE?
7. (a) Suppose you have access to a function FAIRCOIN that returns a single random bit, chosen uniformly and independently from the set {0, 1}, in  $O(1)$  time. Describe and analyze an algorithm RANDOM( $n$ ), which returns an integer chosen uniformly and independently at random from the set  $\{1, 2, \dots, n\}$ .
- (b) Suppose you have access to a function FAIRCOINS( $k$ ) that returns  $k$  random bits (or equivalently, a random integer chosen uniformly and independently from the set  $\{0, 1, \dots, 2^k - 1\}$ ) in  $O(1)$  time, given any non-negative integer  $k$  as input. Describe and analyze an algorithm RANDOM( $n$ ), which returns an integer chosen uniformly and independently at random from the set  $\{1, 2, \dots, n\}$ .

**For each of the remaining problems, you may assume a function RANDOM( $k$ ) that returns, given any positive integer  $k$ , an integer chosen independently and uniformly at random from the set  $\{1, 2, \dots, k\}$ , in  $O(1)$  time. For example, to perform a fair coin flip, one could call RANDOM(2).**

8. Consider the following algorithm for finding the smallest element in an unsorted array:

```

RANDOMMIN( $A[1..n]$ ):
   $min \leftarrow \infty$ 
  for  $i \leftarrow 1$  to  $n$  in random order
    if  $A[i] < min$ 
       $min \leftarrow A[i]$  (*)  

  return  $min$ 

```

- (a) In the worst case, how many times does RANDOMMIN execute line (\*)?
- (b) What is the probability that line (\*) is executed during the  $n$ th iteration of the for loop?
- (c) What is the *exact* expected number of executions of line (\*)?

9. Consider the following randomized algorithm for choosing the largest bolt. Draw a bolt uniformly at random from the set of  $n$  bolts, and draw a nut uniformly at random from the set of  $n$  nuts. If the bolt is smaller than the nut, discard the bolt, draw a new bolt uniformly at random from the unchosen bolts, and repeat. Otherwise, discard the nut, draw a new nut uniformly at random from the unchosen nuts, and repeat. Stop either when every nut has been discarded, or every bolt except the one in your hand has been discarded.

What is the *exact* expected number of nut-bolt tests performed by this algorithm? Prove your answer is correct. [Hint: What is the expected number of unchosen nuts and bolts when the algorithm terminates?]

10. Let  $S$  be a set of  $n$  points in the plane. A point  $p$  in  $S$  is called *Pareto-optimal* if no other point in  $S$  is both above and to the right of  $p$ .

- (a) Describe and analyze a deterministic algorithm that computes the Pareto-optimal points in  $S$  in  $O(n \log n)$  time.
- (b) Suppose each point in  $S$  is chosen independently and uniformly at random from the unit square  $[0, 1] \times [0, 1]$ . What is the *exact* expected number of Pareto-optimal points in  $S$ ?

11. Suppose we want to write an efficient function  $\text{RANDOMPERMUTATION}(n)$  that returns a permutation of the integers  $\langle 1, \dots, n \rangle$  chosen uniformly at random.

- (a) Prove that the following algorithm is *not* correct. [Hint: Consider the case  $n = 3$ .]

```
RANDOMPERMUTATION( $n$ ):
  for  $i \leftarrow 1$  to  $n$ 
     $\pi[i] \leftarrow i$ 
  for  $i \leftarrow 1$  to  $n$ 
    swap  $\pi[i] \leftrightarrow \pi[\text{RANDOM}(n)]$ 
```

- (b) Consider the following implementation of  $\text{RANDOMPERMUTATION}$ .

```
RANDOMPERMUTATION( $n$ ):
  for  $i \leftarrow 1$  to  $n$ 
     $\pi[i] \leftarrow \text{NULL}$ 
  for  $i \leftarrow 1$  to  $n$ 
     $j \leftarrow \text{RANDOM}(n)$ 
    while ( $\pi[j] \neq \text{NULL}$ )
       $j \leftarrow \text{RANDOM}(n)$ 
     $\pi[j] \leftarrow i$ 
  return  $\pi$ 
```

Prove that this algorithm is correct. Analyze its expected runtime.

- (c) Consider the following partial implementation of  $\text{RANDOMPERMUTATION}$ .

```
RANDOMPERMUTATION( $n$ ):
  for  $i \leftarrow 1$  to  $n$ 
     $A[i] \leftarrow \text{RANDOM}(n)$ 
   $\pi \leftarrow \text{SOMEFUNCTION}(A)$ 
  return  $\pi$ 
```

Prove that if the subroutine  $\text{SOMEFUNCTION}$  is deterministic, then this algorithm cannot be correct. [Hint: There is a one-line proof.]

- (d) Describe and analyze an implementation of RANDOMPERMUTATION that runs in expected worst-case time  $O(n)$ .
- (e) Describe and analyze an implementation of RANDOMPERMUTATION that runs in expected worst-case time  $O(n \log n)$ , using fair coin flips (instead of RANDOM) as the only source of randomness.
- \*(f) Consider a correct implementation of RANDOMPERMUTATION( $n$ ) with the following property: whenever it calls RANDOM( $k$ ), the argument  $k$  is at most  $m$ . Prove that this algorithm **always** calls RANDOM at least  $\Omega(\frac{n \log n}{\log m})$  times.
12. A *data stream* is an extremely long sequence of items that you can only read only once, in order. A good example of a data stream is the sequence of packets that pass through a router. Data stream algorithms must process each item in the stream quickly, using very little memory; there is simply too much data to store, and it arrives too quickly for any complex computations. Every data stream algorithm looks roughly like this:

```
DoSomethingInteresting(stream S):
repeat
     $x \leftarrow$  next item in  $S$ 
     $\langle\langle$ do something fast with  $x$  $\rangle\rangle$ 
until  $S$  ends
return  $\langle\langle$ something $\rangle\rangle$ 
```

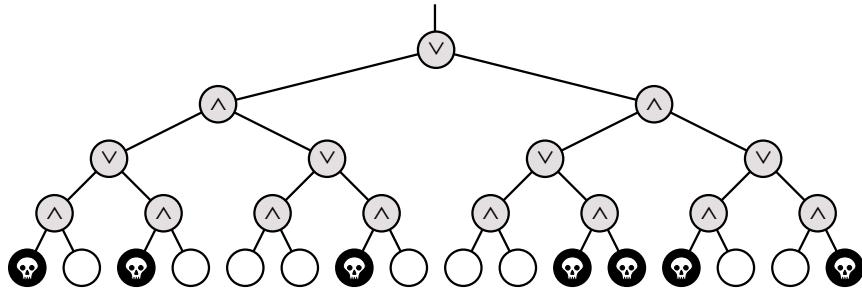
Describe and analyze an algorithm that chooses one element uniformly at random from a data stream, *without knowing the length of the stream in advance*. Your algorithm should spend  $O(1)$  time per stream element and use  $O(1)$  space (not counting the stream itself).

13. The following randomized algorithm, sometimes called ‘one-armed quicksort’, selects the  $r$ th smallest element in an unsorted array  $A[1..n]$ . For example, to find the smallest element, you would call RANDOMSELECT( $A, 1$ ); to find the median element, you would call RANDOMSELECT( $A, \lfloor n/2 \rfloor$ ). Recall from lecture that PARTITION splits the array into three parts by comparing the pivot element  $A[p]$  to every other element of the array, using  $n - 1$  comparisons altogether, and returns the new index of the pivot element. The subroutine RANDOM( $n$ ) returns an integer chosen uniformly at random between 1 and  $n$ , in  $O(1)$  time.

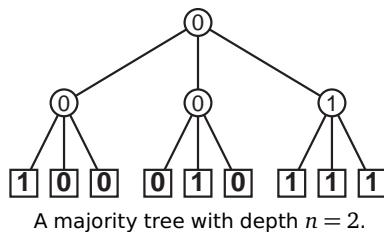
```
RANDOMSELECT( $A[1..n], r$ ):
 $k \leftarrow$  PARTITION( $A[1..n]$ , RANDOM( $n$ ))
if  $r < k$ 
    return RANDOMSELECT( $A[1..k - 1], r$ )
else if  $r > k$ 
    return RANDOMSELECT( $A[k + 1..n], r - k$ )
else
    return  $A[k]$ 
```

- (a) State a recurrence for the expected running time of RANDOMSELECT, as a function of  $n$  and  $r$ .
- (b) What is the *exact* probability that RANDOMSELECT compares the  $i$ th smallest and  $j$ th smallest elements in the input array? The correct answer is a simple function of  $i$ ,  $j$ , and  $r$ . [Hint: Check your answer by trying a few small examples.]

- (c) Show that for any  $n$  and  $r$ , the expected running time of RANDOMSELECT is  $\Theta(n)$ . You can use either the recurrence from part (a) or the probabilities from part (b). For extra credit, find the *exact* expected number of comparisons, as a function of  $n$  and  $r$ .
- (d) What is the expected number of times that RANDOMSELECT calls itself recursively?
14. Let  $M[1..n, 1..n]$  be an  $n \times n$  matrix in which every row and every column is sorted. Such an array is called *totally monotone*. No two elements of  $M$  are equal.
- Describe and analyze an algorithm to solve the following problem in  $O(n)$  time: Given indices  $i, j, i', j'$  as input, compute the number of elements of  $M$  smaller than  $M[i, j]$  and larger than  $M[i', j']$ .
  - Describe and analyze an algorithm to solve the following problem in  $O(n)$  time: Given indices  $i, j, i', j'$  as input, return an element of  $M$  chosen uniformly at random from the elements smaller than  $M[i, j]$  and larger than  $M[i', j']$ . Assume the requested range is always non-empty.
  - Describe and analyze a randomized algorithm to compute the median element of  $M$  in  $O(n \log n)$  expected time.
15. Suppose we have a circular linked list of numbers, implemented as a pair of arrays, one storing the actual numbers and the other storing successor pointers. Specifically, let  $X[1..n]$  be an array of  $n$  distinct real numbers, and let  $N[1..n]$  be an array of indices with the following property: If  $X[i]$  is the largest element of  $X$ , then  $X[N[i]]$  is the smallest element of  $X$ ; otherwise,  $X[N[i]]$  is the smallest element of  $X$  that is larger than  $X[i]$ . For example:
- | $i$    | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  |
|--------|----|----|----|----|----|----|----|----|----|
| $X[i]$ | 83 | 54 | 16 | 31 | 45 | 99 | 78 | 62 | 27 |
| $N[i]$ | 6  | 8  | 9  | 5  | 2  | 3  | 1  | 7  | 4  |
- Describe and analyze a randomized algorithm that determines whether a given number  $x$  appears in the array  $X$  in  $O(\sqrt{n})$  expected time. **Your algorithm may not modify the arrays  $X$  and  $N$ .**
16. Death knocks on your door one cold blustery morning and challenges you to a game. Death knows that you are an algorithms student, so instead of the traditional game of chess, Death presents you with a complete binary tree with  $4^n$  leaves, each colored either black or white. There is a token at the root of the tree. To play the game, you and Death will take turns moving the token from its current node to one of its children. The game will end after  $2n$  moves, when the token lands on a leaf. If the final leaf is black, you die; if it's white, you will live forever. You move first, so Death gets the last turn.
- You can decide whether it's worth playing or not as follows. Imagine that the nodes at even levels (where it's your turn) are OR gates, the nodes at odd levels (where it's Death's turn) are AND gates. Each gate gets its input from its children and passes its output to its parent. White and black stand for TRUE and FALSE. If the output at the top of the tree is TRUE, then you can win and live forever! If the output at the top of the tree is FALSE, you should challenge Death to a game of Twister instead.
- Describe and analyze a deterministic algorithm to determine whether or not you can win. *[Hint: This is easy!]*



- (b) Unfortunately, Death won't give you enough time to look at every node in the tree. Describe a randomized algorithm that determines whether you can win in  $O(3^n)$  expected time. [Hint: Consider the case  $n = 1$ .]
- \*(c) Describe and analyze a randomized algorithm that determines whether you can win in  $O(c^n)$  expected time, for some constant  $c < 3$ . [Hint: You may not need to change your algorithm from part (b) at all!]
17. A *majority tree* is a complete binary tree with depth  $n$ , where every leaf is labeled either 0 or 1. The *value* of a leaf is its label; the *value* of any internal node is the majority of the values of its three children. Consider the problem of computing the value of the root of a majority tree, given the sequence of  $3^n$  leaf labels as input. For example, if  $n = 2$  and the leaves are labeled 1, 0, 0, 0, 1, 0, 1, 1, 1, the root has value 0.



A majority tree with depth  $n = 2$ .

- (a) Prove that *any* deterministic algorithm that computes the value of the root of a majority tree must examine every leaf. [Hint: Consider the special case  $n = 1$ . Recurse.]
- (b) Describe and analyze a randomized algorithm that computes the value of the root in worst-case expected time  $O(c^n)$  for some constant  $c < 3$ . [Hint: Consider the special case  $n = 1$ . Recurse.]

*I thought the following four [rules] would be enough, provided that I made a firm and constant resolution not to fail even once in the observance of them. The first was never to accept anything as true if I had not evident knowledge of its being so. . . . The second, to divide each problem I examined into as many parts as was feasible, and as was requisite for its better solution. The third, to direct my thoughts in an orderly way. . . establishing an order in thought even when the objects had no natural priority one to another. And the last, to make throughout such complete enumerations and such general surveys that I might be sure of leaving nothing out.*

— René Descartes, *Discours de la Méthode* (1637)

*What is luck?*

*Luck is probability taken personally.*

*It is the excitement of bad math.*

— Penn Jillette (2001), quoting Chip Denman (1998)

## 10 Randomized Binary Search Trees

In this lecture, we consider two randomized alternatives to balanced binary search tree structures such as AVL trees, red-black trees, B-trees, or splay trees, which are arguably simpler than any of these deterministic structures.

### 10.1 Treaps

#### 10.1.1 Definitions

A *treap* is a binary tree in which every node has both a *search key* and a *priority*, where the inorder sequence of search keys is sorted and each node's priority is smaller than the priorities of its children.<sup>1</sup> In other words, a treap is simultaneously a binary search tree for the search keys and a (min-)heap for the priorities. In our examples, we will use letters for the search keys and numbers for the priorities.

I'll assume from now on that all the keys and priorities are distinct. Under this assumption, we can easily prove by induction that the structure of a treap is completely determined by the search keys and priorities of its nodes. Since it's a heap, the node  $v$  with highest priority must be the root. Since it's also a binary search tree, any node  $u$  with  $\text{key}(u) < \text{key}(v)$  must be in the left subtree, and any node  $w$  with  $\text{key}(w) > \text{key}(v)$  must be in the right subtree. Finally, since the subtrees are treaps, by induction, their structures are completely determined. The base case is the trivial empty treap.

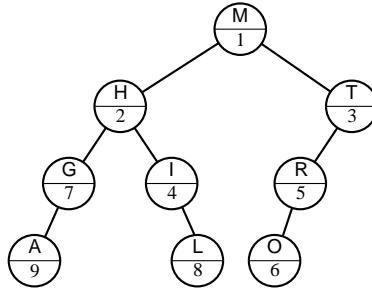
Another way to describe the structure is that a treap is exactly the binary tree that results by inserting the nodes one at a time into an initially empty tree, in order of increasing priority, using the usual insertion algorithm. This is also easy to prove by induction.

A third description interprets the keys and priorities as the coordinates of a set of points in the plane. The root corresponds to a T whose joint lies on the topmost point. The T splits the plane into three parts. The top part is (by definition) empty; the left and right parts are split recursively. This interpretation has some interesting applications in computational geometry, which (unfortunately) we probably won't have time to talk about.

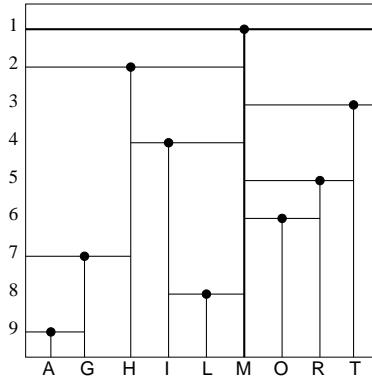
Treaps were first discovered by Jean Vuillemin in 1980, but he called them *Cartesian trees*.<sup>2</sup> The word 'treap' was first used by Edward McCreight around 1980 to describe a slightly different data structure,

<sup>1</sup>Sometimes I hate English. Normally, 'higher priority' means 'more important', but 'first priority' is also more important than 'second priority'. Maybe 'posteriority' would be better; one student suggested 'unimportance'.

<sup>2</sup>J. Vuillemin, A unifying look at data structures. *Commun. ACM* 23:229–239, 1980.



A treap. The top half of each node shows its search key and the bottom half shows its priority.



A geometric interpretation of the same treap.

but he later switched to the more prosaic name *priority search trees*.<sup>3</sup> Treaps were rediscovered and used to build randomized search trees by Cecilia Aragon and Raimund Seidel in 1989.<sup>4</sup> A different kind of randomized binary search tree, which uses random rebalancing instead of random priorities, was later discovered and analyzed by Conrado Martínez and Salvador Roura in 1996.<sup>5</sup>

### 10.1.2 Treap Operations

The search algorithm is the usual one for binary search trees. The time for a successful search is proportional to the depth of the node. The time for an unsuccessful search is proportional to the depth of either its successor or its predecessor.

To insert a new node  $z$ , we start by using the standard binary search tree insertion algorithm to insert it at the bottom of the tree. At the point, the search keys still form a search tree, but the priorities may no longer form a heap. To fix the heap property, as long as  $z$  has smaller priority than its parent, perform a *rotation* at  $z$ , a local operation that decreases the depth of  $z$  by one and increases its parent's depth by one, while maintaining the search tree property. Rotations can be performed in constant time, since they only involve simple pointer manipulation.

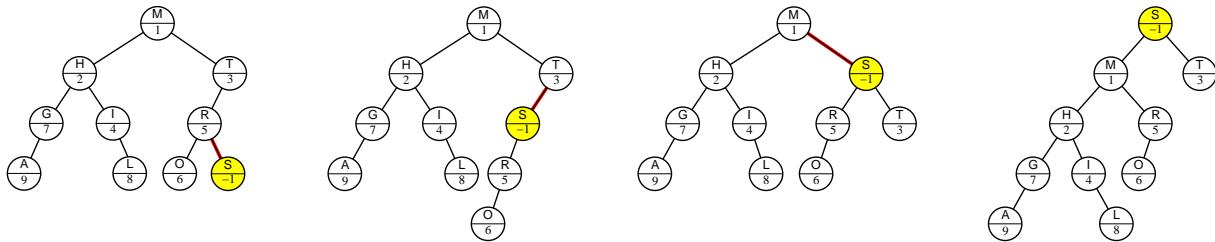
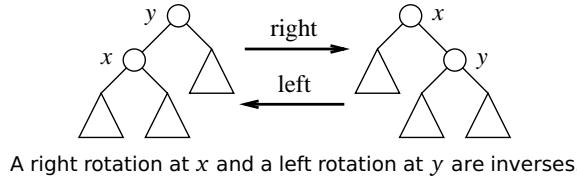
The overall time to insert  $z$  is proportional to the depth of  $z$  before the rotations—we have to walk down the treap to insert  $z$ , and then walk back up the treap doing rotations. Another way to say this is that the time to insert  $z$  is roughly twice the time to perform an unsuccessful search for  $\text{key}(z)$ .

Deleting a node is *exactly* like inserting a node, but in reverse order. Suppose we want to delete node  $z$ . As long as  $z$  is not a leaf, perform a rotation at the child of  $z$  with smaller priority. This moves  $z$

<sup>3</sup>E. M. McCreight. Priority search trees. *SIAM J. Comput.* 14(2):257–276, 1985.

<sup>4</sup>R. Seidel and C. R. Aragon. Randomized search trees. *Algorithmica* 16:464–497, 1996.

<sup>5</sup>C. Martínez and S. Roura. Randomized binary search trees. *J. ACM* 45(2):288–323, 1998. The results in this paper are virtually identical (including the constant factors!) to the corresponding results for treaps, although the analysis techniques are quite different.



down a level and its smaller-priority child up a level. The choice of which child to rotate preserves the heap property everywhere except at  $z$ . When  $z$  becomes a leaf, chop it off.

We sometimes also want to *split* a treap  $T$  into two treaps  $T_<$  and  $T_>$  along some pivot key  $\pi$ , so that all the nodes in  $T_<$  have keys less than  $\pi$  and all the nodes in  $T_>$  have keys bigger than  $\pi$ . A simple way to do this is to insert a new node  $z$  with  $key(z) = \pi$  and  $priority(z) = -\infty$ . After the insertion, the new node is the root of the treap. If we delete the root, the left and right sub-treaps are exactly the trees we want. The time to split at  $\pi$  is roughly twice the time to (unsuccessfully) search for  $\pi$ .

Similarly, we may want to *merge* two treaps  $T_<$  and  $T_>$ , where every node in  $T_<$  has a smaller search key than any node in  $T_>$ , into one super-treap. Merging is just splitting in reverse—create a dummy root whose left sub-treap is  $T_<$  and whose right sub-treap is  $T_>$ , rotate the dummy node down to a leaf, and then cut it off.

The cost of each of these operations is proportional to the depth of some node  $v$  in the treap.

- **Search:** A successful search for key  $k$  takes  $O(depth(v))$  time, where  $v$  is the node with  $key(v) = k$ . For an unsuccessful search, let  $v^-$  be the inorder *predecessor* of  $k$  (the node whose key is just barely smaller than  $k$ ), and let  $v^+$  be the inorder *successor* of  $k$  (the node whose key is just barely larger than  $k$ ). Since the last node examined by the binary search is either  $v^-$  or  $v^+$ , the time for an unsuccessful search is either  $O(depth(v^+))$  or  $O(depth(v^-))$ .
- **Insert/Delete:** Inserting a new node with key  $k$  takes either  $O(depth(v^+))$  time or  $O(depth(v^-))$  time, where  $v^+$  and  $v^-$  are the predecessor and successor of the new node. Deletion is just insertion in reverse.
- **Split/Merge:** Splitting a treap at pivot value  $k$  takes either  $O(depth(v^+))$  time or  $O(depth(v^-))$  time, since it costs the same as inserting a new dummy root with search key  $k$  and priority  $-\infty$ . Merging is just splitting in reverse.

Since the depth of a node in a treap is  $\Theta(n)$  in the worst case, each of these operations has a worst-case running time of  $\Theta(n)$ .

### 10.1.3 Random Priorities

A *randomized treap* is a treap in which the priorities are *independently and uniformly distributed continuous random variables*. That means that whenever we insert a new search key into the treap, we generate a

random real number between (say) 0 and 1 and use that number as the priority of the new node. The only reason we're using real numbers is so that the probability of two nodes having the same priority is zero, since equal priorities make the analysis slightly messier. In practice, we could just choose random integers from a large range, like 0 to  $2^{31} - 1$ , or random floating point numbers. Also, since the priorities are independent, each node is equally likely to have the smallest priority.

The cost of all the operations we discussed—search, insert, delete, split, join—is proportional to the depth of some node in the tree. Here we'll see that the *expected* depth of *any* node is  $O(\log n)$ , which implies that the expected running time for any of those operations is also  $O(\log n)$ .

Let  $x_k$  denote the node with the  $k$ th smallest search key. To analyze the expected depth, we define an indicator variable

$$A_k^i = [x_i \text{ is a proper ancestor of } x_k].$$

(The superscript doesn't mean power in this case; it just a reminder of which node is supposed to be further up in the tree.) Since the depth of  $v$  is just the number of proper ancestors of  $v$ , we have the following identity:

$$\text{depth}(x_k) = \sum_{i=1}^n A_k^i.$$

Now we can express the *expected* depth of a node in terms of these indicator variables as follows.

$$\mathbb{E}[\text{depth}(x_k)] = \sum_{i=1}^n \Pr[A_k^i = 1]$$

(Just as in our analysis of matching nuts and bolts, we're using linearity of expectation and the fact that  $\mathbb{E}[X] = \Pr[X = 1]$  for any indicator variable  $X$ .) So to compute the expected depth of a node, we just have to compute the probability that some node is a proper ancestor of some other node.

Fortunately, we can do this easily once we prove a simple structural lemma. Let  $X(i, k)$  denote either the subset of treap nodes  $\{x_i, x_{i+1}, \dots, x_k\}$  or the subset  $\{x_k, x_{k+1}, \dots, x_i\}$ , depending on whether  $i < k$  or  $i > k$ . The order of the arguments is unimportant; the subsets  $X(i, k)$  and  $X(k, i)$  are identical. The subset  $X(1, n) = X(n, 1)$  contains all  $n$  nodes in the treap.

**Lemma 1.** *For all  $i \neq k$ , node  $x_i$  is a proper ancestor of  $x_k$  if and only if  $x_i$  has the smallest priority among all nodes in  $X(i, k)$ .*

**Proof:** If  $x_i$  is the root, then it is an ancestor of  $x_k$ , and by definition, it has the smallest priority of *any* node in the treap, so it must have the smallest priority in  $X(i, k)$ .

On the other hand, if  $x_k$  is the root, then  $x_i$  is not an ancestor of  $x_k$ , and indeed  $x_i$  does not have the smallest priority in  $X(i, k)$  —  $x_k$  does.

On the gripping hand<sup>6</sup>, suppose some other node  $x_j$  is the root. If  $x_i$  and  $x_k$  are in different subtrees, then either  $i < j < k$  or  $i > j > k$ , so  $x_j \in X(i, k)$ . In this case,  $x_i$  is not an ancestor of  $x_k$ , and indeed  $x_i$  does not have the smallest priority in  $X(i, k)$  —  $x_j$  does.

Finally, if  $x_i$  and  $x_k$  are in the same subtree, the lemma follows inductively (or, if you prefer, recursively), since the subtree is a smaller treap. The empty treap is the trivial base case.  $\square$

Since each node in  $X(i, k)$  is equally likely to have smallest priority, we immediately have the probability we wanted:

---

<sup>6</sup>See Larry Niven and Jerry Pournelle, *The Gripping Hand*, Pocket Books, 1994.

$$\Pr[A_k^i = 1] = \frac{[i \neq k]}{|k - i| + 1} = \begin{cases} \frac{1}{k - i + 1} & \text{if } i < k \\ 0 & \text{if } i = k \\ \frac{1}{i - k + 1} & \text{if } i > k \end{cases}$$

To compute the expected depth of a node  $x_k$ , we just plug this probability into our formula and grind through the algebra.

$$\begin{aligned} \mathbb{E}[depth(x_k)] &= \sum_{i=1}^n \Pr[A_k^i = 1] = \sum_{i=1}^{k-1} \frac{1}{k - i + 1} + \sum_{i=k+1}^n \frac{1}{i - k + 1} \\ &= \sum_{j=2}^k \frac{1}{j} + \sum_{i=2}^{n-k+1} \frac{1}{j} \\ &= H_k - 1 + H_{n-k+1} - 1 \\ &< \ln k + \ln(n - k + 1) - 2 \\ &< 2\ln n - 2. \end{aligned}$$

In conclusion, every search, insertion, deletion, split, and merge operation in an  $n$ -node randomized binary search tree takes  $O(\log n)$  expected time.

Since a treap is exactly the binary tree that results when you insert the keys in order of increasing priority, a randomized treap is the result of inserting the keys in *random* order. So our analysis also automatically gives us the expected depth of any node in a binary tree built by random insertions (without using priorities).

#### 10.1.4 Randomized Quicksort (Again!)

We've already seen two completely different ways of describing randomized quicksort. The first is the familiar recursive one: choose a random pivot, partition, and recurse. The second is a less familiar iterative version: repeatedly choose a new random pivot, partition whatever subset contains it, and continue. But there's a third way to describe randomized quicksort, this time in terms of binary search trees.

RANDOMIZEDQUICKSORT:

$T \leftarrow$  an empty binary search tree  
 insert the keys into  $T$  in *random order*  
 output the inorder sequence of keys in  $T$

Our treap analysis tells us is that this algorithm will run in  $O(n \log n)$  expected time, since each key is inserted in  $O(\log n)$  expected time.

Why is this quicksort? Just like last time, all we've done is rearrange the order of the comparisons. Intuitively, the binary tree is just the recursion tree created by the normal version of quicksort. In the recursive formulation, we compare the initial pivot against everything else and then recurse. In the binary tree formulation, the first "pivot" becomes the root of the tree without any comparisons, but then later as each other key is inserted into the tree, it is compared against the root. Either way, the first pivot chosen is compared with everything else. The partition splits the remaining items into a left subarray and a right subarray; in the binary tree version, these are exactly the items that go into the left

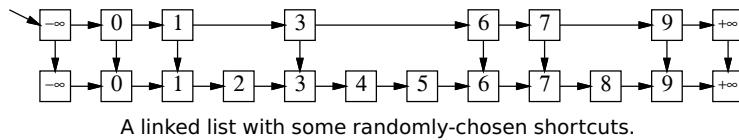
subtree and the right subtree. Since both algorithms define the same two subproblems, by induction, both algorithms perform the same comparisons.

We even saw the probability  $\frac{1}{|k-i|+1}$  before, when we were talking about sorting nuts and bolts with a variant of randomized quicksort. In the more familiar setting of sorting an array of numbers, the probability that randomized quicksort compares the  $i$ th largest and  $k$ th largest elements is exactly  $\frac{2}{|k-i|+1}$ . The binary tree version compares  $x_i$  and  $x_k$  if and only if  $x_i$  is an ancestor of  $x_k$  or vice versa, so the probabilities are exactly the same.

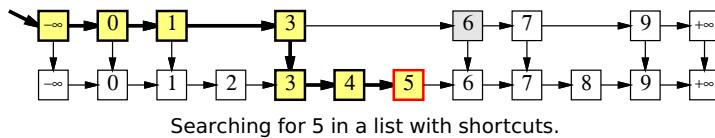
## 10.2 Skip Lists

*Skip lists*, which were first discovered by Bill Pugh in the late 1980's,<sup>7</sup> have many of the usual desirable properties of balanced binary search trees, but their structure is very different.

At a high level, a skip list is just a sorted linked list with some random shortcuts. To do a search in a normal singly-linked list of length  $n$ , we obviously need to look at  $n$  items in the worst case. To speed up this process, we can make a second-level list that contains roughly half the items from the original list. Specifically, for each item in the original list, we duplicate it with probability  $1/2$ . We then string together all the duplicates into a second sorted linked list, and add a pointer from each duplicate back to its original. Just to be safe, we also add sentinel nodes at the beginning and end of both lists.



Now we can find a value  $x$  in this augmented structure using a two-stage algorithm. First, we scan for  $x$  in the shortcut list, starting at the  $-\infty$  sentinel node. If we find  $x$ , we're done. Otherwise, we reach some value bigger than  $x$  and we know that  $x$  is not in the shortcut list. Let  $w$  be the largest item less than  $x$  in the shortcut list. In the second phase, we scan for  $x$  in the original list, starting from  $w$ . Again, if we reach a value bigger than  $x$ , we know that  $x$  is not in the data structure.



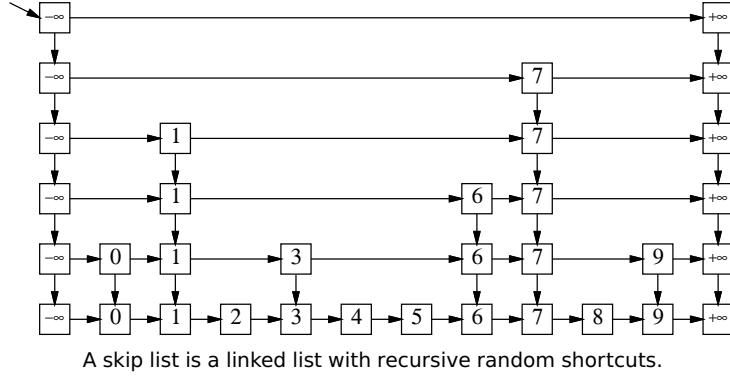
Since each node appears in the shortcut list with probability  $1/2$ , the expected number of nodes examined in the first phase is at most  $n/2$ . Only one of the nodes examined in the second phase has a duplicate. The probability that any node is followed by  $k$  nodes without duplicates is  $2^{-k}$ , so the expected number of nodes examined in the second phase is at most  $1 + \sum_{k \geq 0} 2^{-k} = 2$ . Thus, by adding these random shortcuts, we've reduced the cost of a search from  $n$  to  $n/2 + 2$ , roughly a factor of two in savings.

### 10.2.1 Recursive Random Shortcuts

Now there's an obvious improvement—add shortcuts to the shortcuts, and repeat recursively. That's exactly how skip lists are constructed. For each node in the original list, we flip a coin over and over until we get tails. For each heads, we make a duplicate of the node. The duplicates are stacked up in levels, and the nodes on each level are strung together into sorted linked lists. Each node  $v$  stores a

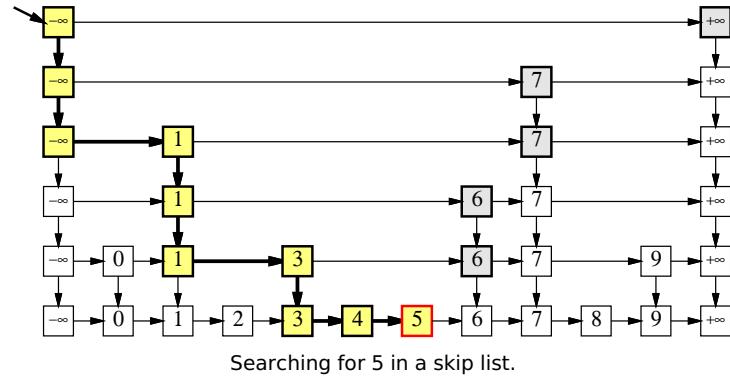
<sup>7</sup>William Pugh. Skip lists: A probabilistic alternative to balanced trees. *Commun. ACM* 33(6):668–676, 1990.

search key ( $key(v)$ ), a pointer to its next lower copy ( $down(v)$ ), and a pointer to the next node in its level ( $right(v)$ ).



The search algorithm for skip lists is very simple. Starting at the leftmost node  $L$  in the highest level, we scan through each level as far as we can without passing the target value  $x$ , and then proceed down to the next level. The search ends when we either reach a node with search key  $x$  or fail to find  $x$  on the lowest level.

```
SKIPLISTFIND( $x, L$ ):
     $v \leftarrow L$ 
    while ( $v \neq \text{NULL}$  and  $key(v) \neq x$ )
        if  $key(right(v)) > x$ 
             $v \leftarrow down(v)$ 
        else
             $v \leftarrow right(v)$ 
    return  $v$ 
```



Intuitively, since each level of the skip lists has about half the number of nodes as the previous level, the total number of levels should be about  $O(\log n)$ . Similarly, each time we add another level of random shortcuts to the skip list, we cut the search time roughly in half, except for a constant overhead, so after  $O(\log n)$  levels, we should have a search time of  $O(\log n)$ . Let's formalize each of these two intuitive observations.

### 10.2.2 Number of Levels

The actual values of the search keys don't affect the skip list analysis, so let's assume the keys are the integers 1 through  $n$ . Let  $L(x)$  be the number of levels of the skip list that contain some search key  $x$ , not

counting the bottom level. Each new copy of  $x$  is created with probability  $1/2$  from the previous level, essentially by flipping a coin. We can compute the expected value of  $L(x)$  recursively—with probability  $1/2$ , we flip tails and  $L(x) = 0$ ; and with probability  $1/2$ , we flip heads, increase  $L(x)$  by one, and recurse:

$$E[L(x)] = \frac{1}{2} \cdot 0 + \frac{1}{2}(1 + E[L(x)])$$

Solving this equation gives us  $E[L(x)] = 1$ .

In order to analyze the expected worst-case cost of a search, however, we need a bound on the *number of levels*  $L = \max_x L(x)$ . Unfortunately, we can't compute the average of a maximum the way we would compute the average of a sum. Instead, we will derive a stronger result, showing that the depth is  $O(\log n)$  **with high probability**. ‘High probability’ is a technical term that means the probability is at least  $1 - 1/n^c$  for some constant  $c \geq 1$ ; the hidden constant in the  $O(\log n)$  bound could depend on  $c$ .

In order for a search key  $x$  to appear on level  $\ell$ , it must have flipped  $\ell$  heads in a row when it was inserted, so  $\Pr[L(x) \geq \ell] = 2^{-\ell}$ . The skip list has at least  $\ell$  levels if and only if  $L(x) \geq \ell$  for at least one of the  $n$  search keys.

$$\Pr[L \geq \ell] = \Pr[(L(1) \geq \ell) \vee (L(2) \geq \ell) \vee \dots \vee (L(n) \geq \ell)]$$

Using the *union bound* —  $\Pr[A \vee B] \leq \Pr[A] + \Pr[B]$  for any random events  $A$  and  $B$  — we can simplify this as follows:

$$\Pr[L \geq \ell] \leq \sum_{x=1}^n \Pr[L(x) \geq \ell] = n \cdot \Pr[L(x) \geq \ell] = \frac{n}{2^\ell}.$$

When  $\ell \leq \lg n$ , this bound is trivial. However, for any constant  $c > 1$ , we have a strong upper bound

$$\Pr[L \geq c \lg n] \leq \frac{1}{n^{c-1}}.$$

We conclude that **with high probability, a skip list has  $O(\log n)$  levels**.

This high-probability bound indirectly implies a bound on the *expected* number of levels. Some simple algebra gives us the following alternate definition for expectation:

$$E[L] = \sum_{\ell \geq 0} \ell \cdot \Pr[L = \ell] = \sum_{\ell \geq 1} \Pr[L \geq \ell]$$

Clearly, if  $\ell < \ell'$ , then  $\Pr[L(x) \geq \ell] > \Pr[L(x) \geq \ell']$ . So we can derive an upper bound on the expected number of levels as follows:

$$\begin{aligned} E[L(x)] &= \sum_{\ell \geq 1} \Pr[L \geq \ell] \\ &= \sum_{\ell=1}^{\lg n} \Pr[L \geq \ell] + \sum_{\ell \geq \lg n+1} \Pr[L \geq \ell] \\ &\leq \sum_{\ell=1}^{\lg n} 1 + \sum_{\ell \geq \lg n+1} \frac{n}{2^\ell} \\ &= \lg n + \sum_{i \geq 1} \frac{1}{2^i} \quad [i = \ell - \lg n] \\ &= \lg n + 2 \end{aligned}$$

So in expectation, a skip list has *at most two* more levels than an ideal version where each level contains exactly half the nodes of the next level below.

### 10.2.3 Logarithmic Search Time

It's a little easier to analyze the cost of a search if we imagine running the algorithm backwards. `SKIPLISTFIND` takes the output from `SKIPLISTFIND` as input and traces back through the data structure to the upper left corner. Skip lists don't really have up and left pointers, but we'll pretend that they do so we don't have to write '(v)nowob → v' or '(v)righting → v'.<sup>8</sup>

```
SKIPLISTFIND(v):
    while ( $v \neq L$ )
        if  $up(v)$  exists
             $v \leftarrow up(v)$ 
        else
             $v \leftarrow left(v)$ 
```

Now for *every* node  $v$  in the skip list,  $up(v)$  exists with probability  $1/2$ . So for purposes of analysis, `SKIPLISTFIND` is equivalent to the following algorithm:

```
FLIPWALK(v):
    while ( $v \neq L$ )
        if COINFLIP = HEADS
             $v \leftarrow up(v)$ 
        else
             $v \leftarrow left(v)$ 
```

Obviously, the expected number of heads is exactly the same as the expected number of TAILS. Thus, the expected running time of this algorithm is twice the expected number of upward jumps. Since we already know that the number of upward jumps is  $O(\log n)$  with high probability, we can conclude that the worst-case search time is  $O(\log n)$  with high probability (and therefore in expectation).

## Exercises

1. (a) The *left spine* of a binary tree is a path starting at the root and following only left-child pointers down to a leaf. What is the expected number of nodes in the left spine of an  $n$ -node treap?  
(b) What is the expected number of leaves in an  $n$ -node treap? [Hint: What is the probability that in an  $n$ -node treap, the node with  $k$ th smallest search key is a leaf?]  
(c) Prove that the expected number of proper descendants of any node in a treap is exactly equal to the expected depth of that node.
2. Consider an  $n$ -node treap  $T$ . As in the lecture notes, we identify nodes in  $T$  by the ranks of their search keys. Thus, 'node 5' means the node with the 5th smallest search key. Let  $i, j, k$  be integers such that  $1 \leq i \leq j \leq k \leq n$ .
  - (a) What is the *exact* probability that node  $j$  is a common ancestor of node  $i$  and node  $k$ ?
  - (b) What is the *exact* expected length of the unique path from node  $i$  to node  $k$  in  $T$ ?

---

Learned about VNCI while all this notes using monitor-writing, but not because he wanted to keep his discoveries secret.<sup>8</sup>  
He just had really bad timing in his right hand!

3. Recall that a *priority search tree* is a binary tree in which every node has both a *search key* and a *priority*, arranged so that the tree is simultaneously a binary search tree for the keys and a min-heap for the priorities. A **heater** is a priority search tree in which the *priorities* are given by the user, and the *search keys* are distributed uniformly and independently at random in the real interval  $[0, 1]$ . Intuitively, a heater is a sort of dual treap.

The following problems consider an  $n$ -node heater  $T$  whose node priorities are the integers from 1 to  $n$ . We identify nodes in  $T$  by their priorities; thus, ‘node 5’ means the node in  $T$  with priority 5. The min-heap property implies that node 1 is the root of  $T$ . Finally, let  $i$  and  $j$  be integers with  $1 \leq i < j \leq n$ .

- (a) Prove that in a random permutation of the  $(i + 1)$ -element set  $\{1, 2, \dots, i, j\}$ , elements  $i$  and  $j$  are adjacent with probability  $2/(i + 1)$ .
- (b) Prove that node  $i$  is an ancestor of node  $j$  with probability  $2/(i + 1)$ . [Hint: Use part (a)!]
- (c) What is the probability that node  $i$  is a descendant of node  $j$ ? [Hint: Don’t use part (a)!]
- (d) What is the *exact* expected depth of node  $j$ ?
- (e) Describe and analyze an algorithm to insert a new item into a heater. Express the expected running time of the algorithm in terms of the rank of the newly inserted item.
- (f) Describe an algorithm to delete the minimum-priority item (the root) from an  $n$ -node heater. What is the expected running time of your algorithm?

- \*4. In the usual theoretical presentation of treaps, the priorities are random real numbers chosen uniformly from the interval  $[0, 1]$ , but in practice, computers only have access to random *bits*. This problem asks you to analyze a modification of treaps that takes this limitation into account.

Suppose the priority of a node  $v$  is abstractly represented as an infinite sequence  $\pi_v[1 .. \infty]$  of random bits, which is interpreted as the rational number

$$\text{priority}(v) = \sum_{i=1}^{\infty} \pi_v[i] \cdot 2^{-i}.$$

However, only a finite number  $\ell_v$  of these bits are actually known at any given time. When a node  $v$  is first created, *none* of the priority bits are known:  $\ell_v = 0$ . We generate (or ‘reveal’) new random bits only when they are necessary to compare priorities. The following algorithm compares the priorities of any two nodes in  $O(1)$  expected time:

```

LARGERPRIORITY( $v, w$ ):
  for  $i \leftarrow 1$  to  $\infty$ 
    if  $i > \ell_v$ 
       $\ell_v \leftarrow i$ ;  $\pi_v[i] \leftarrow \text{RANDOMBIT}$ 
    if  $i > \ell_w$ 
       $\ell_w \leftarrow i$ ;  $\pi_w[i] \leftarrow \text{RANDOMBIT}$ 
    if  $\pi_v[i] > \pi_w[i]$ 
      return  $v$ 
    else if  $\pi_v[i] < \pi_w[i]$ 
      return  $w$ 
  
```

Suppose we insert  $n$  items one at a time into an initially empty treap. Let  $L = \sum_v \ell_v$  denote the total number of random bits generated by calls to LARGERPRIORITY during these insertions.

- (a) Prove that  $E[L] = \Theta(n)$ .
- (b) Prove that  $E[\ell_v] = \Theta(1)$  for any node  $v$ . [Hint: This is equivalent to part (a). Why?]
- (c) Prove that  $E[\ell_{\text{root}}] = \Theta(\log n)$ . [Hint: Why doesn't this contradict part (b)?]
- \*5. Any skip list  $\mathcal{L}$  can be transformed into a binary search tree  $T(\mathcal{L})$  as follows. The root of  $T(\mathcal{L})$  is the leftmost node on the highest non-empty level of  $\mathcal{L}$ ; the left and right subtrees are constructed recursively from the nodes to the left and to the right of the root. Let's call the resulting tree  $T(\mathcal{L})$  a *skip list tree*.
- (a) Show that any search in  $T(\mathcal{L})$  is no more expensive than the corresponding search in  $\mathcal{L}$ . (Searching in  $T(\mathcal{L})$  could be *considerably* cheaper—why?)
- (b) Describe an algorithm to insert a new search key into a skip list tree in  $O(\log n)$  expected time. Inserting key  $x$  into  $T(\mathcal{L})$  should produce *exactly* the same tree as inserting  $x$  into  $\mathcal{L}$  and then transforming  $\mathcal{L}$  into a tree. [Hint: You will need to maintain some additional information in the tree nodes.]
- (c) Describe an algorithm to delete a search key from a skip list tree in  $O(\log n)$  expected time. Again, deleting key  $x$  from  $T(\mathcal{L})$  should produce *exactly* the same tree as deleting  $x$  from  $\mathcal{L}$  and then transforming  $\mathcal{L}$  into a tree.
6. A *meldable priority queue* stores a set of keys from some totally-ordered universe (such as the integers) and supports the following operations:
- MAKEQUEUE: Return a new priority queue containing the empty set.
  - FINDMIN( $Q$ ): Return the smallest element of  $Q$  (if any).
  - DELETEMIN( $Q$ ): Remove the smallest element in  $Q$  (if any).
  - INSERT( $Q, x$ ): Insert element  $x$  into  $Q$ , if it is not already there.
  - DECREASEKEY( $Q, x, y$ ): Replace an element  $x \in Q$  with a smaller key  $y$ . (If  $y > x$ , the operation fails.) The input is a pointer directly to the node in  $Q$  containing  $x$ .
  - DELETE( $Q, x$ ): Delete the element  $x \in Q$ . The input is a pointer directly to the node in  $Q$  containing  $x$ .
  - MELD( $Q_1, Q_2$ ): Return a new priority queue containing all the elements of  $Q_1$  and  $Q_2$ ; this operation destroys  $Q_1$  and  $Q_2$ .

A simple way to implement such a data structure is to use a heap-ordered binary tree, where each node stores a key, along with pointers to its parent and two children. MELD can be implemented using the following randomized algorithm:

```

MELD( $Q_1, Q_2$ ):
  if  $Q_1$  is empty return  $Q_2$ 
  if  $Q_2$  is empty return  $Q_1$ 
  if  $\text{key}(Q_1) > \text{key}(Q_2)$ 
    swap  $Q_1 \leftrightarrow Q_2$ 
  with probability 1/2
     $\text{left}(Q_1) \leftarrow \text{MELD}(\text{left}(Q_1), Q_2)$ 
  else
     $\text{right}(Q_1) \leftarrow \text{MELD}(\text{right}(Q_1), Q_2)$ 
  return  $Q_1$ 

```

- (a) Prove that for *any* heap-ordered binary trees  $Q_1$  and  $Q_2$  (not just those constructed by the operations listed above), the expected running time of  $\text{MELD}(Q_1, Q_2)$  is  $O(\log n)$ , where  $n = |Q_1| + |Q_2|$ . [*Hint: How long is a random root-to-leaf path in an  $n$ -node binary tree if each left/right choice is made with equal probability?*]
- (b) Prove that  $\text{MELD}(Q_1, Q_2)$  runs in  $O(\log n)$  time with high probability.
- (c) Show that each of the other meldable priority queue operations can be implemented with at most one call to  $\text{MELD}$  and  $O(1)$  additional time. (This implies that every operation takes  $O(\log n)$  time with high probability.)

*If you hold a cat by the tail  
you learn things you cannot learn any other way.*  
— Mark Twain

## \*11 Tail Inequalities

The simple recursive structure of skip lists made it relatively easy to derive an upper bound on the expected *worst-case* search time, by way of a stronger high-probability upper bound on the worst-case search time. We can prove similar results for treaps, but because of the more complex recursive structure, we need slightly more sophisticated probabilistic tools. These tools are usually called *tail inequalities*; intuitively, they bound the probability that a random variable with a bell-shaped distribution takes a value in the *tails* of the distribution, far away from the mean.

### 11.1 Markov's Inequality

Perhaps the simplest tail inequality was named after the Russian mathematician Andrey Markov; however, in strict accordance with Stigler's Law of Eponymy, it first appeared in the works of Markov's probability teacher, Pafnuty Chebyshev.<sup>1</sup>

**Markov's Inequality.** Let  $X$  be a non-negative integer random variable. For any  $t > 0$ , we have  $\Pr[X \geq t] \leq \mathbb{E}[X]/t$ .

**Proof:** The inequality follows from the definition of expectation by simple algebraic manipulation.

$$\begin{aligned}
 \mathbb{E}[X] &= \sum_{k=0}^{\infty} k \cdot \Pr[X = k] && [\text{definition of } \mathbb{E}[X]] \\
 &= \sum_{k=0}^{\infty} \Pr[X \geq k] && [\text{algebra}] \\
 &\geq \sum_{k=0}^{t-1} \Pr[X \geq k] && [\text{since } t < \infty] \\
 &\geq \sum_{k=0}^{t-1} \Pr[X \geq t] && [\text{since } k < t] \\
 &= t \cdot \Pr[X \geq t] && [\text{algebra}] \quad \square
 \end{aligned}$$

Unfortunately, the bounds that Markov's inequality implies (at least directly) are often very weak, even useless. (For example, Markov's inequality implies that with high probability, every node in an  $n$ -node treap has depth  $O(n^2 \log n)$ . Well, duh!) To get stronger bounds, we need to exploit some additional structure in our random variables.

---

<sup>1</sup>The closely related tail bound traditionally called Chebyshev's inequality was actually discovered by the French statistician Irénée-Jules Bienaymé, a friend and colleague of Chebyshev's.

## 11.2 Independence

A set of random variables  $X_1, X_2, \dots, X_n$  are said to be *mutually independent* if and only if

$$\Pr \left[ \bigwedge_{i=1}^n (X_i = x_i) \right] = \prod_{i=1}^n \Pr[X_i = x_i]$$

for all possible values  $x_1, x_2, \dots, x_n$ . For examples, different flips of the same fair coin are mutually independent, but the number of heads and the number of tails in a sequence of  $n$  coin flips are not independent (since they must add to  $n$ ). Mutual independence of the  $X_i$ 's implies that the expectation of the product of the  $X_i$ 's is equal to the product of the expectations:

$$\mathbb{E} \left[ \prod_{i=1}^n X_i \right] = \prod_{i=1}^n \mathbb{E}[X_i].$$

Moreover, if  $X_1, X_2, \dots, X_n$  are independent, then for any function  $f$ , the random variables  $f(X_1), f(X_2), \dots, f(X_n)$  are also mutually independent.

## 11.3 Chernoff Bounds

Suppose  $X = \sum_{i=1}^n X_i$  is the sum of  $n$  mutually independent random *indicator* variables  $X_i$ . For each  $i$ , let  $p_i = \Pr[X_i = 1]$ , and let  $\mu = \mathbb{E}[X] = \sum_i \mathbb{E}[X_i] = \sum_i p_i$ .

**Chernoff Bound (Upper Tail).**  $\Pr[X > (1 + \delta)\mu] < \left( \frac{e^\delta}{(1 + \delta)^{1+\delta}} \right)^\mu$  for any  $\delta > 0$ .

**Proof:** The proof is fairly long, but it replies on just a few basic components: a clever substitution, Markov's inequality, the independence of the  $X_i$ 's, The World's Most Useful Inequality  $e^x > 1 + x$ , a tiny bit of calculus, and lots of high-school algebra.

We start by introducing a variable  $t$ , whose role will become clear shortly.

$$\Pr[X > (1 + \delta)\mu] = \Pr[e^{tX} > e^{t(1+\delta)\mu}]$$

To cut down on the superscripts, I'll usually write  $\exp(x)$  instead of  $e^x$  in the rest of the proof. Now apply Markov's inequality to the right side of this equation:

$$\Pr[X > (1 + \delta)\mu] < \frac{\mathbb{E}[\exp(tX)]}{\exp(t(1 + \delta)\mu)}.$$

We can simplify the expectation on the right using the fact that the terms  $X_i$  are independent.

$$\mathbb{E}[\exp(tX)] = \mathbb{E} \left[ \exp \left( t \sum_i X_i \right) \right] = \mathbb{E} \left[ \prod_i \exp(tX_i) \right] = \prod_i \mathbb{E}[\exp(tX_i)]$$

We can bound the individual expectations  $\mathbb{E}[\exp(tX_i)]$  using The World's Most Useful Inequality:

$$\mathbb{E}[\exp(tX_i)] = p_i e^t + (1 - p_i) = 1 + (e^t - 1)p_i < \exp((e^t - 1)p_i)$$

This inequality gives us a simple upper bound for  $\mathbb{E}[e^{tX}]$ :

$$\mathbb{E}[\exp(tX)] < \prod_i \exp((e^t - 1)p_i) < \exp \left( \sum_i (e^t - 1)p_i \right) = \exp((e^t - 1)\mu)$$

Substituting this back into our original fraction from Markov's inequality, we obtain

$$\Pr[X > (1 + \delta)\mu] < \frac{\mathbb{E}[\exp(tX)]}{\exp(t(1 + \delta)\mu)} < \frac{\exp((e^t - 1)\mu)}{\exp(t(1 + \delta)\mu)} = (\exp(e^t - 1 - t(1 + \delta)))^\mu$$

Notice that this last inequality holds for *all* possible values of  $t$ . To obtain the final tail bound, we will choose  $t$  to make this bound as small as possible. To minimize  $e^t - 1 - t - t\delta$ , we take its derivative with respect to  $t$  and set it to zero:

$$\frac{d}{dt}(e^t - 1 - t(1 + \delta)) = e^t - 1 - \delta = 0.$$

(And you thought calculus would never be useful!) This equation has just one solution  $t = \ln(1 + \delta)$ . Plugging this back into our bound gives us

$$\Pr[X > (1 + \delta)\mu] < (\exp(\delta - (1 + \delta)\ln(1 + \delta)))^\mu = \left(\frac{e^\delta}{(1 + \delta)^{1+\delta}}\right)^\mu$$

And we're done! □

This form of the Chernoff bound can be a bit clumsy to use. A more complicated argument gives us the bound

$$\boxed{\Pr[X > (1 + \delta)\mu] < e^{-\mu\delta^2/3} \text{ for any } 0 < \delta < 1.}$$

A similar argument gives us an inequality bounding the probability that  $X$  is significantly *smaller* than its expected value:

**Chernoff Bound (Lower Tail).**  $\boxed{\Pr[X < (1 - \delta)\mu] < \left(\frac{e^{-\delta}}{(1 - \delta)^{1-\delta}}\right)^\mu < e^{-\mu\delta^2/2} \text{ for any } \delta > 0.}$

## 11.4 Back to Treaps

In our analysis of randomized treaps, we defined the indicator variable  $A_k^i$  to have the value 1 if and only if the node with the  $i$ th smallest key ('node  $i$ ') was a proper ancestor of the node with the  $k$ th smallest key ('node  $k$ '). We argued that

$$\Pr[A_k^i = 1] = \frac{[i \neq k]}{|k - i| + 1},$$

and from this we concluded that the expected depth of node  $k$  is

$$\mathbb{E}[\text{depth}(k)] = \sum_{i=1}^n \Pr[A_k^i = 1] = H_k + H_{n-k} - 2 < 2 \ln n.$$

To prove a worst-case expected bound on the depth of the tree, we need to argue that the *maximum* depth of any node is small. Chernoff bounds make this argument easy, once we establish that the relevant indicator variables are mutually independent.

**Lemma 1.** *For any index  $k$ , the  $k-1$  random variables  $A_k^i$  with  $i < k$  are mutually independent. Similarly, for any index  $k$ , the  $n-k$  random variables  $A_k^i$  with  $i > k$  are mutually independent.*

**Proof:** To simplify the notation, we explicitly consider only the case  $k = 1$ , although the argument generalizes easily to other values of  $k$ . Fix  $n - 1$  arbitrary indicator values  $x_2, x_3, \dots, x_n$ . We prove the lemma by induction on  $n$ , with the vacuous base case  $n = 1$ . The definition of conditional probability gives us

$$\begin{aligned} \Pr \left[ \bigwedge_{i=2}^n (A_1^i = x_i) \right] &= \Pr \left[ \bigwedge_{i=2}^{n-1} (A_k^i = x_i) \wedge A_1^n = x_n \right] \\ &= \Pr \left[ \bigwedge_{i=2}^{n-1} (A_k^i = x_i) \mid A_1^n = x_n \right] \cdot \Pr [A_1^n = x_n] \end{aligned}$$

Now recall that  $A_1^n = 1$  if and only if node  $n$  has the smallest priority, and the other  $n - 2$  indicator variables  $A_1^i$  depend only on the order of the priorities of nodes 1 through  $n - 1$ . There are exactly  $(n - 1)!$  permutations of the  $n$  priorities in which the  $n$ th priority is smallest, and each of these permutations is equally likely. Thus,

$$\Pr \left[ \bigwedge_{i=2}^{n-1} (A_k^i = x_i) \mid A_1^n = x_n \right] = \Pr \left[ \bigwedge_{i=2}^{n-1} (A_k^i = x_i) \right]$$

The inductive hypothesis implies that the variables  $A_1^2, \dots, A_1^{n-1}$  are mutually independent, so

$$\Pr \left[ \bigwedge_{i=2}^{n-1} (A_k^i = x_i) \right] = \prod_{i=2}^{n-1} \Pr [A_1^i = x_i].$$

We conclude that

$$\Pr \left[ \bigwedge_{i=2}^n (A_1^i = x_i) \right] = \Pr [A_1^n = x_n] \cdot \prod_{i=2}^{n-1} \Pr [A_1^i = x_i] = \prod_{i=1}^{n-1} \Pr [A_1^i = x_i],$$

or in other words, that the indicator variables are mutually independent.  $\square$

**Theorem 2.** *The depth of a randomized treap with  $n$  nodes is  $O(\log n)$  with high probability.*

**Proof:** First let's bound the probability that the depth of node  $k$  is at most  $8 \ln n$ . There's nothing special about the constant 8 here; I'm being generous to make the analysis easier.

The depth is a sum of  $n$  indicator variables  $A_k^i$ , as  $i$  ranges from 1 to  $n$ . Our Observation allows us to partition these variables into two mutually independent subsets. Let  $d_{<}(k) = \sum_{i < k} A_k^i$  and  $d_{>}(k) = \sum_{i < k} A_k^i$ , so that  $\text{depth}(k) = d_{<}(k) + d_{>}(k)$ . If  $\text{depth}(k) > 8 \ln n$ , then either  $d_{<}(k) > 4 \ln n$  or  $d_{>}(k) > 4 \ln n$ .

Chernoff's inequality, with  $\mu = \mathbb{E}[d_{<}(k)] = H_k - 1 < \ln n$  and  $\delta = 3$ , bounds the probability that  $d_{<}(k) > 4 \ln n$  as follows.

$$\Pr[d_{<}(k) > 4 \ln n] < \Pr[d_{<}(k) > 4\mu] < \left( \frac{e^3}{4^4} \right)^\mu < \left( \frac{e^3}{4^4} \right)^{\ln n} = n^{\ln(e^3/4^4)} = n^{3-4\ln 4} < \frac{1}{n^2}.$$

(The last step uses the fact that  $4 \ln 4 \approx 5.54518 > 5$ .) The same analysis implies that  $\Pr[d_{>}(k) > 4 \ln n] < 1/n^2$ . These inequalities imply the crude bound  $\Pr[\text{depth}(k) > 4 \ln n] < 2/n^2$ .

Now consider the probability that the treap has depth greater than  $10 \ln n$ . Even though the distributions of different nodes' depths are *not* independent, we can conservatively bound the probability of failure as follows:

$$\Pr \left[ \max_k \text{depth}(k) > 8 \ln n \right] = \Pr \left[ \bigwedge_{k=1}^n (\text{depth}(k) > 8 \ln n) \right] \leq \sum_{k=1}^n \Pr[\text{depth}(k) > 8 \ln n] < \frac{2}{n}.$$

This argument implies more generally that for any constant  $c$ , the depth of the treap is greater than  $c \ln n$  with probability at most  $2/n^{c \ln c - c}$ . We can make the failure probability an arbitrarily small polynomial by choosing  $c$  appropriately.  $\square$

This lemma implies that any search, insertion, deletion, or merge operation on an  $n$ -node treap requires  $O(\log n)$  time with high probability. In particular, the expected *worst-case* time for each of these operations is  $O(\log n)$ .

## Exercises

1. Prove that for any integer  $k$  such that  $1 < k < n$ , the  $n - 1$  indicator variables  $A_k^i$  with  $i \neq k$  are *not* mutually independent. [*Hint: Consider the case  $n = 3$ .*]
2. Recall from Exercise 1 in the previous note that the expected number of descendants of any node in a treap is  $O(\log n)$ . Why doesn't the Chernoff-bound argument for depth imply that, with high probability, *every* node in a treap has  $O(\log n)$  descendants? The conclusion is clearly bogus—Every treap has a node with  $n$  descendants!—but what's the hole in the argument?
3. A *heater* is a sort of dual treap, in which the priorities of the nodes are given, but their search keys are generated independently and uniformly from the unit interval  $[0, 1]$ . You can assume all priorities and keys are distinct.
  - (a) Prove that for any  $r$ , the node with the  $r$ th smallest *priority* has expected depth  $O(\log r)$ .
  - (b) Prove that an  $n$ -node heater has depth  $O(\log n)$  with high probability.
  - (c) Describe algorithms to perform the operations `INSERT` and `DELETEMIN` in a heater. What are the expected worst-case running times of your algorithms? In particular, can you express the expected running time of `INSERT` in terms of the priority rank of the newly inserted item?

**Calvin:** There! I finished our secret code!  
**Hobbes:** Let's see.  
**Calvin:** I assigned each letter a totally random number, so the code will be hard to crack. For letter "A", you write 3,004,572,688. "B" is 28,731,569½.  
**Hobbes:** That's a good code all right.  
**Calvin:** Now we just commit this to memory.  
**Calvin:** Did you finish your map of our neighborhood?  
**Hobbes:** Not yet. How many bricks does the front walk have?

— Bill Watterson, “Calvin and Hobbes” (August 23, 1990)

## 12 Hash Tables

### 12.1 Introduction

A *hash table* is a data structure for storing a set of items, so that we can quickly determine whether an item is or is not in the set. The basic idea is to pick a *hash function*  $h$  that maps every possible item  $x$  to a small integer  $h(x)$ . Then we store  $x$  in slot  $h(x)$  in an array. The array is the hash table.

Let's be a little more specific. We want to store a set of  $n$  items. Each item is an element of some finite<sup>1</sup> set  $\mathcal{U}$  called the *universe*; we use  $u$  to denote the size of the universe, which is just the number of items in  $\mathcal{U}$ . A hash table is an array  $T[1..m]$ , where  $m$  is another positive integer, which we call the *table size*. Typically,  $m$  is much smaller than  $u$ . A *hash function* is any function of the form

$$h: \mathcal{U} \rightarrow \{0, 1, \dots, m - 1\},$$

mapping each possible item in  $\mathcal{U}$  to a slot in the hash table. We say that an item  $x$  *hashes* to the slot  $T[h(x)]$ .

Of course, if  $u = m$ , then we can always just use the trivial hash function  $h(x) = x$ . In other words, use the item itself as the index into the table. This is called a *direct access table*, or more commonly, an *array*. In most applications, though, the universe of possible keys is orders of magnitude too large for this approach to be practical. Even when it is possible to allocate enough memory, we usually need to store only a small fraction of the universe. Rather than wasting lots of space, we should make  $m$  roughly equal to  $n$ , the number of items in the set we want to maintain.

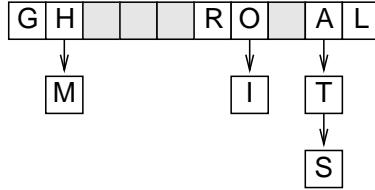
What we'd like is for every item in our set to hash to a different position in the array. Unfortunately, unless  $m = u$ , this is too much to hope for, so we have to deal with *collisions*. We say that two items  $x$  and  $y$  *collide* if they have the same hash value:  $h(x) = h(y)$ . Since we obviously can't store two items in the same slot of an array, we need to describe some methods for *resolving* collisions. The two most common methods are called *chaining* and *open addressing*.

### 12.2 Chaining

In a *chained* hash table, each entry  $T[i]$  is not just a single item, but rather (a pointer to) a linked list of all the items that hash to  $T[i]$ . Let  $\ell(x)$  denote the length of the list  $T[h(x)]$ . To see if an item  $x$  is in the hash table, we scan the entire list  $T[h(x)]$ . The worst-case time required to search for  $x$  is  $O(1)$  to

<sup>1</sup>This finiteness assumption is necessary for several of the technical details to work out, but can be ignored in practice. To hash elements from an infinite universe (for example, the positive integers), pretend that the universe is actually finite but very very large. In fact, in *real* practice, the universe actually *is* finite but very very large. For example, on most modern computers, there are only  $2^{64}$  integers (unless you use a big integer package like GMP, in which case the number of integers is closer to  $2^{2^{32}}$ .)

compute  $h(x)$  plus  $O(1)$  for every element in  $T[h(x)]$ , or  $O(1 + \ell(x))$  overall. Inserting and deleting  $x$  also take  $O(1 + \ell(x))$  time.



A chained hash table with load factor 1.

In the worst case, every item would be hashed to the same value, so we'd get just one long list of  $n$  items. In principle, for any deterministic hashing scheme, a malicious adversary can always present a set of items with exactly this property. In order to defeat such malicious behavior, we'd like to use a hash function that is as random as possible. Choosing a truly random hash function is completely impractical, but there are several heuristics for producing hash functions that behave randomly, or at least close to randomly on real data. Thus, we will analyze the performance as though our hash function were truly random. More formally, we make the following assumption.

**Simple uniform hashing assumption:** If  $x \neq y$  then  $\Pr[h(x) = h(y)] = 1/m$ .

In the next section, I'll describe a small set of functions with the property that a random hash function in this set satisfies the simple uniform hashing assumption. Most actual implementations of hash tables use *deterministic* hash functions. These clearly violate the uniform hashing assumption—the collision probability is either 0 or 1, depending on the pair of items! Nevertheless, it is common practice to adopt the uniform hashing assumption as a convenient fiction for purposes of analysis.

Let's compute the expected value of  $\ell(x)$  under this assumption; this will immediately imply a bound on the expected time to search for an item  $x$ . To be concrete, let's suppose that  $x$  is not already stored in the hash table. For all items  $x$  and  $y$ , we define the indicator variable

$$C_{x,y} = [h(x) = h(y)].$$

(In case you've forgotten the bracket notation,  $C_{x,y} = 1$  if  $h(x) = h(y)$  and  $C_{x,y} = 0$  if  $h(x) \neq h(y)$ .) Since the length of  $T[h(x)]$  is precisely equal to the number of items that collide with  $x$ , we have

$$\ell(x) = \sum_{y \in T} C_{x,y}.$$

We can rewrite the simple uniform hashing assumption as follows:

$$x \neq y \implies E[C_{x,y}] = \Pr[C_{x,y} = 1] = \frac{1}{m}.$$

Now we just have to grind through the definitions.

$$E[\ell(x)] = \sum_{y \in T} E[C_{x,y}] = \sum_{y \in T} \frac{1}{m} = \frac{n}{m}$$

We call this fraction  $n/m$  the *load factor* of the hash table. Since the load factor shows up everywhere, we will give it its own symbol  $\alpha$ .

$$\alpha := \frac{n}{m}$$

Our analysis implies that the expected time for an unsuccessful search in a chained hash table is  $\Theta(1 + \alpha)$ . As long as the number of items  $n$  is only a constant factor bigger than the table size  $m$ , the search time is a constant. A similar analysis gives the same expected time bound (with a slightly smaller constant) for a successful search.

Obviously, linked lists are not the only data structure we could use to store the chains; any data structure that can store a set of items will work. For example, if the universe  $\mathcal{U}$  has a total ordering, we can store each chain in a balanced binary search tree. This reduces the expected time for any search to  $O(1 + \log \ell(x))$ , and under the simple uniform hashing assumption, the expected time for any search is  $O(1 + \log \alpha)$ .

Another natural possibility is to work recursively! Specifically, for each  $T[i]$ , we maintain a hash table  $T_i$  containing all the items with hash value  $i$ . Collisions in those secondary tables are resolved recursively, by storing secondary overflow lists in tertiary hash tables, and so on. The resulting data structure is a tree of hash tables, whose leaves correspond to items that (at some level of the tree) are hashed without any collisions. If every hash table in this tree has size  $m$ , then the expected time for any search is  $O(\log_m n)$ . In particular, if we set  $m = \sqrt{n}$ , the expected time for any search is *constant*. On the other hand, there is no inherent reason to use the same hash table size everywhere; after all, hash tables deeper in the tree are storing fewer items.

**Caveat Lector!**<sup>2</sup> The preceding analysis does *not* imply bounds on the expected worst-case search time is constant. The expected worst-case search time is  $O(1 + L)$ , where  $L = \max_x \ell(x)$ . Under the uniform hashing assumption, the maximum list size  $L$  is *very* likely to grow faster than any constant, unless the load factor  $\alpha$  is *significantly* smaller than 1. For example,  $E[L] = \Theta(\log n / \log \log n)$  when  $\alpha = 1$ . We've stumbled on a powerful but counterintuitive fact about probability: When several individual items are distributed independently and uniformly at random, the resulting distribution is *not* uniform in the traditional sense! Later in this lecture, I'll describe how to achieve constant expected worst-case search time using secondary hash tables.

### 12.3 Universal Hashing

Now I'll describe a method to generate random hash functions that satisfy the simple uniform hashing assumption. We say that a set  $\mathcal{H}$  of hash function is *universal* if it satisfies the following property: For any items  $x \neq y$ , if a hash function  $h$  is chosen *uniformly at random* from the set  $\mathcal{H}$ , then  $\Pr[h(x) = h(y)] = 1/m$ . Note that this probability holds for *any* items  $x$  and  $y$ ; the randomness is entirely in choosing a hash function from the set  $\mathcal{H}$ .

To simplify the following discussion, I'll assume that the universe  $\mathcal{U}$  contains exactly  $m^2$  items, each represented as a pair  $(x, x')$  of integers between 0 and  $m - 1$ . (Think of the items as two-digit numbers in base  $m$ .) I will also assume that  $m$  is a prime number.

For any integers  $0 \leq a, b \leq m - 1$ , define the function  $h_{a,b} : \mathcal{U} \rightarrow \{0, 1, \dots, m - 1\}$  as follows:

$$h_{a,b}(x, x') = (ax + bx') \bmod m.$$

Then the set

$$\mathcal{H} = \{h_{a,b} \mid 0 \leq a, b \leq m - 1\}$$

of all such functions is universal. To prove this, we need to show that for any pair of distinct items  $(x, x') \neq (y, y')$ , exactly  $m$  of the  $m^2$  functions in  $\mathcal{H}$  cause a collision.

---

<sup>2</sup>Latin for “Beware of cannibals.”

Choose two items  $(x, x') \neq (y, y')$ , and assume without loss of generality<sup>3</sup> that  $x \neq y$ . A function  $h_{a,b} \in \mathcal{H}$  causes a collision between  $(x, x')$  and  $(y, y')$  if and only if

$$\begin{aligned} h_{a,b}(x, x') &= h_{a,b}(y, y') \\ (ax + bx') \bmod m &= (ay + by') \bmod m \\ ax + bx' &\equiv ay + by' \pmod{m} \\ a(x - y) &\equiv b(y' - x') \pmod{m} \\ a &\equiv \frac{b(y' - x')}{x - y} \pmod{m}. \end{aligned}$$

In the last step, we are using the fact that  $m$  is prime and  $x - y \neq 0$ , which implies that  $x - y$  has a unique multiplicative inverse modulo  $m$ . (For example, the multiplicative inverse of 12 modulo 17 is 10, since  $12 \cdot 10 = 120 \equiv 1 \pmod{17}$ .) For each possible value of  $b$ , the last identity defines a *unique* value of  $a$  such that  $h_{a,b}$  causes a collision. Since there are  $m$  possible values for  $b$ , there are exactly  $m$  hash functions  $h_{a,b}$  that cause a collision, which is exactly what we needed to prove.

Thus, if we want to achieve the constant expected time bounds described earlier, we should choose a random element of  $\mathcal{H}$  as our hash function, by generating two numbers  $a$  and  $b$  uniformly at random between 0 and  $m - 1$ . This is *precisely* the same as choosing a element of  $\mathcal{U}$  uniformly at random.

One perhaps undesirable ‘feature’ of this construction is that we have a small chance of choosing the trivial hash function  $h_{0,0}$ , which maps everything to 0. So in practice, if we happen to pick  $a = b = 0$ , we should reject that choice and pick new random numbers. By taking  $h_{0,0}$  out of consideration, we reduce the probability of a collision from  $1/m$  to  $(m - 1)/(m^2 - 1) = 1/(m + 1)$ . In other words, the set  $\mathcal{H} \setminus \{h_{0,0}\}$  is slightly better than universal.

This construction can be easily generalized to larger universes. Suppose  $u = m^r$  for some constant  $r$ , so that each element  $x \in \mathcal{U}$  can be represented by a vector  $(x_0, x_1, \dots, x_{r-1})$  of integers between 0 and  $m - 1$ . (Think of  $x$  as an  $r$ -digit number written in base  $m$ .) Then for each vector  $a = (a_0, a_1, \dots, a_{r-1})$ , define the corresponding hash function  $h_a$  as follows:

$$h_a(x) = (a_0 x_0 + a_1 x_1 + \dots + a_{r-1} x_{r-1}) \bmod m.$$

Then the set of all  $m^r$  such functions is universal.

#### \*12.4 High Probability Bounds: Balls and Bins

Although any particular search in a chained hash tables requires only constant expected time, but what about the *worst* search time? Under a stronger version<sup>4</sup> of the uniform hashing assumption, this is equivalent to the following more abstract problem. Suppose we toss  $n$  balls independently and uniformly at random into one of  $n$  bins. Can we say anything about the number of balls in the fullest bin?

**Lemma 1.** *If  $n$  balls are thrown independently and uniformly into  $n$  bins, then with high probability, the fullest bin contains  $O(\log n / \log \log n)$  balls.*

---

<sup>3</sup>‘Without loss of generality’ is a phrase that appears (perhaps too) often in combinatorial proofs. What it means is that we are considering one of many possible cases, but once we see the proof for one case, the proofs for all the other cases are obvious thanks to some inherent symmetry. For this proof, we are not explicitly considering what happens when  $x = y$  and  $x' \neq y'$ .

<sup>4</sup>The simple uniform hashing assumption requires only *pairwise* independence, but the following analysis requires *full* independence.

**Proof:** Let  $X_j$  denote the number of balls in bin  $j$ , and let  $\hat{X} = \max_j X_j$  be the maximum number of balls in any bin. Clearly,  $E[X_j] = 1$  for all  $j$ .

Now consider the probability that bin  $j$  contains exactly  $k$  balls. There are  $\binom{n}{k}$  choices for those  $k$  balls; each chosen ball has probability  $1/n$  of landing in bin  $j$ ; and each of the remaining balls has probability  $1 - 1/n$  of missing bin  $j$ . Thus,

$$\begin{aligned}\Pr[X_j = k] &= \binom{n}{k} \left(\frac{1}{n}\right)^k \left(1 - \frac{1}{n}\right)^{n-k} \\ &\leq \frac{n^k}{k!} \left(\frac{1}{n}\right)^k \left(1 - \frac{1}{n}\right)^{n-k} \\ &< \frac{1}{k!}\end{aligned}$$

This bound shrinks super-exponentially as  $k$  increases, so we can very crudely bound the probability that bin 1 contains *at least*  $k$  balls as follows:

$$\Pr[X_j \geq k] < \frac{n}{k!}$$

To prove the high-probability bound, we need to choose a value for  $k$  such that  $n/k! \approx 1/n^c$  for some constant  $c$ . Taking logs of both sides of the desired approximation and applying Stirling's approximation, we find

$$\begin{aligned}\ln k! &\approx k \ln k - k \approx (c+1) \ln n \\ \iff k &\approx \frac{(c+1) \ln n}{\ln k - 1} \\ &\approx \frac{(c+1) \ln n}{\ln \frac{(c+1) \ln n}{\ln k - 1} - 1} \\ &= \frac{(c+1) \ln n}{\ln \ln n + \ln(c+1) - \ln(\ln k - 1) - 1} \\ &\approx \frac{(c+1) \ln n}{\ln \ln n}.\end{aligned}$$

We have shown (modulo some algebraic hand-waving that is easy but tedious to clean up) that

$$\Pr\left[X_j \geq \frac{(c+1) \ln n}{\ln \ln n}\right] < \frac{1}{n^c}.$$

This probability bound holds for every bin  $j$ . Thus, by the union bound, we conclude that

$$\begin{aligned}\Pr\left[\max_j X_j > \frac{(c+1) \ln n}{\ln \ln n}\right] &= \Pr\left[X_j > \frac{(c+1) \ln n}{\ln \ln n} \text{ for all } j\right] \\ &\leq \sum_{j=1}^n \Pr\left[X_j > \frac{(c+1) \ln n}{\ln \ln n}\right] \\ &< \frac{1}{n^{c-1}}.\end{aligned}\quad \square$$

A somewhat more complicated argument implies that if we throw  $n$  balls randomly into  $n$  bins, then with high probability, the most popular bin contains at least  $\Omega(\log n / \log \log n)$  balls.

However, if we make the hash table large enough, we can expect every ball to land in its own bin. Suppose there are  $m$  bins. Let  $C_{ij}$  be the indicator variable that equals 1 if and only if  $i \neq j$  and ball  $i$  and ball  $j$  land in the same bin, and let  $C = \sum_{i < j} C_{ij}$  be the total number of pairwise collisions. Since the balls are thrown uniformly at random, the probability of a collision is exactly  $1/m$ , so  $E[C] = \binom{n}{2}/m$ . In particular, if  $m = n^2$ , the expected number of collisions is less than  $1/2$ .

To get a high probability bound, let  $X_j$  denote the number of balls in bin  $j$ , as in the previous proof. We can easily bound the probability that bin  $j$  is empty, by taking the two most significant terms in a binomial expansion:

$$\Pr[X_j = 0] = \left(1 - \frac{1}{m}\right)^n = \sum_{i=1}^n \binom{n}{i} \left(\frac{-1}{m}\right)^i = 1 - \frac{n}{m} + \Theta\left(\frac{n^2}{m^2}\right) > 1 - \frac{n}{m}$$

We can similarly bound the probability that bin  $j$  contains exactly one ball:

$$\Pr[X_j = 1] = n \cdot \frac{1}{m} \left(1 - \frac{1}{m}\right)^{n-1} = \frac{n}{m} \left(1 - \frac{n-1}{m} + \Theta\left(\frac{n^2}{m^2}\right)\right) > \frac{n}{m} - \frac{n(n-1)}{m^2}$$

It follows immediately that  $\Pr[X_j > 1] < n(n-1)/m^2$ . The union bound now implies that  $\Pr[\hat{X} > 1] < n(n-1)/m$ . If we set  $m = n^{2+\varepsilon}$  for any constant  $\varepsilon > 0$ , then the probability that no bin contains more than one ball is at least  $1 - 1/n^\varepsilon$ .

**Lemma 2.** *For any  $\varepsilon > 0$ , if  $n$  balls are thrown independently and uniformly into  $n^{2+\varepsilon}$  bins, then with high probability, no bin contains more than one ball.*

## 12.5 Perfect Hashing

So far we are faced with two alternatives. If we use a small hash table, to keep the space usage down, the resulting worst-case expected search time is  $\Theta(\log n / \log \log n)$  with high probability, which is not much better than a binary search tree. On the other hand, we can get constant worst-case search time, at least in expectation, by using a table of quadratic size, but that seems unduly wasteful.

Fortunately, there is a fairly simple way to combine these two ideas to get a data structure of linear expected size, whose expected worst-case search time is constant. At the top level, we use a hash table of size  $n$ , but instead of linked lists, we use secondary hash tables to resolve collisions. Specifically, the  $j$ th secondary hash table has size  $n_j^2$ , where  $n_j$  is the number of items whose primary hash value is  $j$ . The expected worst-case search time in any secondary hash table is  $O(1)$ , by our earlier analysis.

Although this data structure needs significantly more memory for each secondary structure, the overall increase in space is insignificant, at least in expectation.

**Lemma 3.** *The simple uniform hashing assumption implies  $E[n_j^2] < 2$ .*

**Proof:** Let  $X_{ij}$  be the indicator variable that equals 1 if item  $i$  hashes to slot  $j$  in the primary hash table. Linearity of expectation implies that

$$E[n_j^2] = E\left[\sum_{i=1}^n \sum_{k=1}^n X_{ij} X_{kj}\right] = \sum_{i=1}^n \sum_{k=1}^n E[X_{ij} X_{kj}] = \sum_{i=1}^n E[X_{ij}^2] + 2 \sum_{i=1}^n \sum_{k=i+1}^n E[X_{ij} X_{kj}].$$

Because  $X_{ij}$  is an indicator variable, we have  $X_{ij}^2 = X_{ij}$ , which implies that  $E[X_{ij}^2] = E[X_{ij}] = 1/n$  by the uniform hashing assumption. The uniform hashing assumption also implies that  $X_{ij}$  and  $X_{kj}$  are

independent whenever  $i \neq k$ , so  $E[X_{ij}X_{kj}] = E[X_{ij}]E[X_{kj}] = 1/n^2$ . Thus,

$$E[n_j^2] = \sum_{i=1}^n \frac{1}{n} + 2 \sum_{i=1}^n \sum_{k=i+1}^n \frac{1}{n^2} = 1 + 2 \binom{n}{2} \frac{1}{n^2} = 2 - \frac{1}{n}. \quad \square$$

This lemma implies that the expected size of our two-level hash table is  $O(n)$ . By our earlier analysis, the expected worst-case search time is  $O(1)$ .

## 12.6 Open Addressing

Another method used to resolve collisions in hash tables is called *open addressing*. Here, rather than building secondary data structures, we resolve collisions by looking elsewhere in the table. Specifically, we have a sequence of hash functions  $\langle h_0, h_1, h_2, \dots, h_{m-1} \rangle$ , such that for any item  $x$ , the *probe sequence*  $\langle h_0(x), h_1(x), \dots, h_{m-1}(x) \rangle$  is a permutation of  $\langle 0, 1, 2, \dots, m-1 \rangle$ . In other words, different hash functions in the sequence always map  $x$  to different locations in the hash table.

We search for  $x$  using the following algorithm, which returns the array index  $i$  if  $T[i] = x$ , ‘absent’ if  $x$  is not in the table but there is an empty slot, and ‘full’ if  $x$  is not in the table and there no no empty slots.

```
OPENADDRESSSEARCH( $x$ ):
  for  $i \leftarrow 0$  to  $m - 1$ 
    if  $T[h_i(x)] = x$ 
      return  $h_i(x)$ 
    else if  $T[h_i(x)] = \emptyset$ 
      return ‘absent’
  return ‘full’
```

The algorithm for inserting a new item into the table is similar; only the second-to-last line is changed to  $T[h_i(x)] \leftarrow x$ . Notice that for an open-addressed hash table, the load factor is never bigger than 1.

Just as with chaining, we’d like to pretend that the sequence of hash values is random. For purposes of analysis, there is a stronger uniform hashing assumption that gives us constant expected search and insertion time.

### Strong uniform hashing assumption:

For any item  $x$ , the probe sequence  $\langle h_0(x), h_1(x), \dots, h_{m-1}(x) \rangle$  is equally likely to be any permutation of the set  $\{0, 1, 2, \dots, m-1\}$ .

Let’s compute the expected time for an unsuccessful search using this stronger assumption. Suppose there are currently  $n$  elements in the hash table. Our strong uniform hashing assumption has two important consequences:

- The initial hash value  $h_0(x)$  is equally likely to be any integer in the set  $\{0, 1, 2, \dots, m-1\}$ .
- If we ignore the first probe, the remaining probe sequence  $\langle h_1(x), h_2(x), \dots, h_{m-1}(x) \rangle$  is equally likely to be any permutation of the smaller set  $\{0, 1, 2, \dots, m-1\} \setminus \{h_0(x)\}$ .

The first sentence implies that the probability that  $T[h_0(x)]$  is occupied is exactly  $n/m$ . The second sentence implies that if  $T[h_0(x)]$  is occupied, *our search algorithm recursively searches the rest of the hash table!* Since the algorithm will never again probe  $T[h_0(x)]$ , for purposes of analysis, we might as well

pretend that slot in the table no longer exists. Thus, we get the following recurrence for the expected number of probes, as a function of  $m$  and  $n$ :

$$\mathbb{E}[T(m, n)] = 1 + \frac{n}{m} \mathbb{E}[T(m - 1, n - 1)].$$

The trivial base case is  $T(m, 0) = 1$ ; if there's nothing in the hash table, the first probe always hits an empty slot. We can now easily prove by induction that

$\text{EMPHE}[T(m, n)] \leq m/(m - n)$ :

$$\begin{aligned} \mathbb{E}[T(m, n)] &= 1 + \frac{n}{m} \mathbb{E}[T(m - 1, n - 1)] \\ &\leq 1 + \frac{n}{m} \cdot \frac{m - 1}{m - n} && [\text{induction hypothesis}] \\ &< 1 + \frac{n}{m} \cdot \frac{m}{m - n} && [m - 1 < m] \\ &= \frac{m}{m - n} && [\text{algebra}] \end{aligned}$$

Rewriting this in terms of the load factor  $\alpha = n/m$ , we get  $\mathbb{E}[T(m, n)] \leq 1/(1 - \alpha)$ . In other words, the expected time for an unsuccessful search is  $O(1)$ , unless the hash table is almost completely full.

In practice, however, we can't generate truly random probe sequences, so we use one of the following heuristics:

- **Linear probing:** We use a single hash function  $h(x)$ , and define  $h_i(x) = (h(x) + i) \bmod m$ . This is nice and simple, but collisions tend to make items in the table clump together badly, so this is not really a good idea.
- **Quadratic probing:** We use a single hash function  $h(x)$ , and define  $h_i(x) = (h(x) + i^2) \bmod m$ . Unfortunately, for certain values of  $m$ , the sequence of hash values  $\langle h_i(x) \rangle$  does not hit every possible slot in the table; we can avoid this problem by making  $m$  a prime number. (That's often a good idea anyway.) Although quadratic probing does not suffer from the same clumping problems as linear probing, it does have a weaker clustering problem: If two items have the same initial hash value, their entire probe sequences will be the same.
- **Double hashing:** We use two hash functions  $h(x)$  and  $h'(x)$ , and define  $h_i$  as follows:

$$h_i(x) = (h(x) + i \cdot h'(x)) \bmod m$$

To guarantee that this can hit every slot in the table, the *stride* function  $h'(x)$  and the table size  $m$  must be relatively prime. We can guarantee this by making  $m$  prime, but a simpler solution is to make  $m$  a power of 2 and choose a stride function that is always odd. Double hashing avoids the clustering problems of linear and quadratic probing. In fact, the actual performance of double hashing is almost the same as predicted by the uniform hashing assumption, at least when  $m$  is large and the component hash functions  $h$  and  $h'$  are sufficiently random. This is the method of choice!<sup>5</sup>

## 12.7 Cuckoo Hashing

To be written!

<sup>5</sup>...unless your hash tables are really huge, in which case linear probing has *far* better cache behavior, especially when the load factor is small.

## Exercises

1. Suppose we are using an *open-addressed* hash table of size  $m$  to store  $n$  items, where  $n \leq m/2$ . Assume that the strong uniform hashing assumption holds. For any  $i$ , let  $X_i$  denote the number of probes required for the  $i$ th insertion into the table, and let  $X = \max_i X_i$  denote the length of the longest probe sequence.
  - (a) Prove that  $\Pr[X_i > k] \leq 1/2^k$  for all  $i$  and  $k$ .
  - (b) Prove that  $\Pr[X_i > 2\lg n] \leq 1/n^2$  for all  $i$ .
  - (c) Prove that  $\Pr[X > 2\lg n] \leq 1/n$ .
  - (d) Prove that  $E[X] = O(\lg n)$ .
2. Your boss wants you to find a *perfect* hash function for mapping a known set of  $n$  items into a table of size  $m$ . A hash function is *perfect* if there are *no* collisions; each of the  $n$  items is mapped to a different slot in the hash table. Of course, this requires that  $m \geq n$ . (This is a different definition of perfect hashing than the one considered in the lecture notes.) After cursing your algorithms instructor for not teaching you about perfect hashing, you decide to try something simple: repeatedly pick random hash functions until you find one that happens to be perfect. A random hash function  $h$  satisfies two properties:
  - $\Pr[h(x) = h(y)] = 1/m$  for any pair of items  $x \neq y$ .
  - $\Pr[h(x) = i] = 1/m$  for any item  $x$  and any integer  $1 \leq i \leq m$ .
  - (a) Suppose you pick a random hash function  $h$ . What is the *exact* expected number of collisions, as a function of  $n$  (the number of items) and  $m$  (the size of the table)? Don't worry about how to resolve collisions; just count them.
  - (b) What is the *exact* probability that a random hash function is perfect?
  - (c) What is the *exact* expected number of different random hash functions you have to test before you find a perfect hash function?
  - (d) What is the *exact* probability that none of the first  $N$  random hash functions you try is perfect?
  - (e) How many random hash functions do you have to test to find a perfect hash function *with high probability*?
- ★3. Recall that  $F_k$  denotes the  $k$ th Fibonacci number:  $F_0 = 0$ ,  $F_1 = 1$ , and  $F_k = F_{k-1} + F_{k-2}$  for all  $k \geq 2$ . Suppose we are building a hash table of size  $m = F_k$  using the hash function

$$h(x) = (F_{k-1} \cdot x) \bmod F_k$$

Prove that if the consecutive integers  $0, 1, 2, \dots, F_k - 1$  are inserted in order into an initially empty table, each integer will be hashed into the largest contiguous empty interval in the table. In particular, show that there are no collisions.

**Jaques:** But, for the seventh cause; how did you find the quarrel on the seventh cause?

**Touchstone:** Upon a lie seven times removed:-bear your body more seeming, Audrey:-as thus, sir. I did dislike the cut of a certain courtier's beard: he sent me word, if I said his beard was not cut well, he was in the mind it was: this is called the Retort Courteous. If I sent him word again 'it was not well cut,' he would send me word, he cut it to please himself: this is called the Quip Modest. If again 'it was not well cut,' he disabled my judgment: this is called the Reply Churlish. If again 'it was not well cut,' he would answer, I spake not true: this is called the Reproof Valiant. If again 'it was not well cut,' he would say I lied: this is called the Counter-cheque Quarrelsome: and so to the Lie Circumstantial and the Lie Direct.

**Jaques:** And how oft did you say his beard was not well cut?

**Touchstone:** I durst go no further than the Lie Circumstantial, nor he durst not give me the Lie Direct; and so we measured swords and parted.

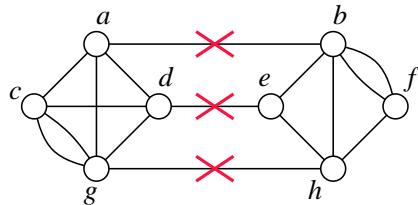
— William Shakespeare, *As You Like It*, Act V, Scene 4 (1600)

## 13 Randomized Minimum Cut

### 13.1 Setting Up the Problem

This lecture considers a problem that arises in robust network design. Suppose we have a connected multigraph<sup>1</sup>  $G$  representing a communications network like the UIUC telephone system, the internet, or Al-Qaeda. In order to disrupt the network, an enemy agent plans to remove some of the edges in this multigraph (by cutting wires, placing police at strategic drop-off points, or paying street urchins to 'lose' messages) to separate it into multiple components. Since his country is currently having an economic crisis, the agent wants to remove as few edges as possible to accomplish this task.

More formally, a *cut* partitions the nodes of  $G$  into two nonempty subsets. The *size* of the cut is the number of *crossing edges*, which have one endpoint in each subset. Finally, a *minimum cut* in  $G$  is a cut with the smallest number of crossing edges. The same graph may have several minimum cuts.



A multigraph whose minimum cut has three edges.

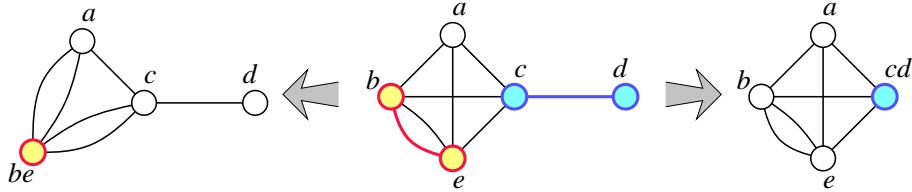
This problem has a long history. The classical deterministic algorithms for this problem rely on *network flow* techniques, which are discussed in another lecture. The fastest such algorithms (that we will discuss) run in  $O(n^3)$  time and are fairly complex; we will see some of these later in the semester. Here I'll describe a relatively simple randomized algorithm discovered by David Karger when he was a Ph.D. student.<sup>2</sup>

Karger's algorithm uses a primitive operation called *collapsing an edge*. Suppose  $u$  and  $v$  are vertices that are connected by an edge in some multigraph  $G$ . To collapse the edge  $\{u, v\}$ , we create a new node called  $uv$ , replace any edge of the form  $\{u, w\}$  or  $\{v, w\}$  with a new edge  $\{uv, w\}$ , and then delete the

<sup>1</sup>A multigraph allows multiple edges between the same pair of nodes. Everything in this lecture could be rephrased in terms of simple graphs where every edge has a non-negative weight, but this would make the algorithms and analysis slightly more complicated.

<sup>2</sup>David R. Karger\*. Random sampling in cut, flow, and network design problems. Proc. 25th STOC, 648–657, 1994.

original vertices  $u$  and  $v$ . Equivalently, collapsing the edge shrinks the edge down to nothing, pulling the two endpoints together. The new collapsed graph is denoted  $G/\{u, v\}$ . We don't allow self-loops in our multigraphs; if there are multiple edges between  $u$  and  $v$ , collapsing any one of them deletes them all.



A graph  $G$  and two collapsed graphs  $G/\{b, e\}$  and  $G/\{c, d\}$ .

Any edge in an  $n$ -vertex graph can be collapsed in  $O(n)$  time, assuming the graph is represented as an adjacency list; I'll leave the precise implementation details as an easy exercise.

The correctness of our algorithms will eventually boil down the following simple observation: For any cut in  $G/\{u, v\}$ , there is cut in  $G$  with exactly the same number of crossing edges. In fact, in some sense, the 'same' edges form the cut in both graphs. The converse is not necessarily true, however. For example, in the picture above, the original graph  $G$  has a cut of size 1, but the collapsed graph  $G/\{c, d\}$  does not.

This simple observation has two immediate but important consequences. First, collapsing an edge cannot decrease the minimum cut size. More importantly, collapsing an edge increases the minimum cut size if and only if that edge is part of *every* minimum cut.

## 13.2 Blindly Guessing

Let's start with an algorithm that tries to guess the minimum cut by randomly collapsing edges until the graph has only two vertices left.

```

GUESSMINCUT( $G$ ):
for  $i \leftarrow n$  downto 2
    pick a random edge  $e$  in  $G$ 
     $G \leftarrow G/e$ 
return the only cut in  $G$ 

```

Because each collapse requires  $O(n)$  time, this algorithm runs in  $O(n^2)$  time. Our earlier observations imply that as long as we never collapse an edge that lies in every minimum cut, our algorithm will actually guess correctly. But how likely is that?

Suppose  $G$  has only one minimum cut—if it actually has more than one, just pick your favorite—and this cut has size  $k$ . Every vertex of  $G$  must lie on at least  $k$  edges; otherwise, we could separate that vertex from the rest of the graph with an even smaller cut. Thus, the number of incident vertex-edge pairs is at least  $kn$ . Since every edge is incident to exactly two vertices,  $G$  must have at least  $kn/2$  edges. That implies that if we pick an edge in  $G$  uniformly at random, the probability of picking an edge in the minimum cut is at most  $2/n$ . In other words, the probability that we don't screw up on the very first step is at least  $1 - 2/n$ .

Once we've collapsed the first random edge, the rest of the algorithm proceeds recursively (with independent random choices) on the remaining  $(n - 1)$ -node graph. So the overall probability  $P(n)$  that GUESSMINCUT returns the true minimum cut is given by the following recurrence:

$$P(n) \geq \frac{n-2}{n} \cdot P(n-1).$$

The base case for this recurrence is  $P(2) = 1$ . We can immediately expand this recurrence into a product, most of whose factors cancel out immediately.

$$P(n) \geq \prod_{i=3}^n \frac{i-2}{i} = \frac{\prod_{i=3}^n (i-2)}{\prod_{i=3}^n i} = \frac{\prod_{i=1}^{n-2} i}{\prod_{i=3}^n i} = \boxed{\frac{2}{n(n-1)}}$$

### 13.3 Blindly Guessing Over and Over

That's not very good. Fortunately, there's a simple method for increasing our chances of finding the minimum cut: run the guessing algorithm many times and return the smallest guess. Randomized algorithms folks like to call this idea *amplification*.

```
KARGERMINCUT( $G$ ):
   $mink \leftarrow \infty$ 
  for  $i \leftarrow 1$  to  $N$ 
     $X \leftarrow \text{GUESSMINCUT}(G)$ 
    if  $|X| < mink$ 
       $mink \leftarrow |X|$ 
       $minX \leftarrow X$ 
  return  $minX$ 
```

Both the running time and the probability of success will depend on the number of iterations  $N$ , which we haven't specified yet.

First let's figure out the probability that KARGERMINCUT returns the actual minimum cut. The only way for the algorithm to return the wrong answer is if GUESSMINCUT fails  $N$  times in a row. Since each guess is independent, our probability of success is at least

$$1 - \left(1 - \frac{2}{n(n-1)}\right)^N \leq 1 - e^{-2N/n(n-1)},$$

by The World's Most Useful Inequality  $1 + x \leq e^x$ . By making  $N$  larger, we can make this probability arbitrarily close to 1, but never equal to 1. In particular, if we set  $N = c \binom{n}{2} \ln n$  for some constant  $c$ , then KARGERMINCUT is correct with probability at least

$$1 - e^{-c \ln n} = 1 - \frac{1}{n^c}.$$

When the failure probability is a polynomial fraction, we say that the algorithm is correct *with high probability*. Thus, KARGERMINCUT computes the minimum cut of any  $n$ -node graph in  $O(n^4 \log n)$  time.

If we make the number of iterations even larger, say  $N = n^2(n-1)/2$ , the success probability becomes  $1 - e^{-n}$ . When the failure probability is exponentially small like this, we say that the algorithm is correct with *very high probability*. In practice, very high probability is usually overkill; high probability is enough. (Remember, there is a small but non-zero probability that your computer will transform itself into a kitten before your program is finished.)

### 13.4 Not-So-Blindly Guessing

The  $O(n^4 \log n)$  running time is actually comparable to some of the simpler flow-based algorithms, but it's nothing to get excited about. But we can improve our guessing algorithm, and thus decrease the

number of iterations in the outer loop, by observing that *as the graph shrinks, the probability of collapsing an edge in the minimum cut increases*. At first the probability is quite small, only  $2/n$ , but near the end of execution, when the graph has only three vertices, we have a  $2/3$  chance of screwing up!

A simple technique for working around this increasing probability of error was developed by David Karger and Cliff Stein.<sup>3</sup> Their idea is to group the first several random collapses a ‘safe’ phase, so that the cumulative probability of screwing up is small—less than  $1/2$ , say—and a ‘dangerous’ phase, which is much more likely to screw up.

The safe phase shrinks the graph from  $n$  nodes to  $n/\sqrt{2} + 1$  nodes, using a sequence of  $n - n/\sqrt{2} - 1$  random collapses. Following our earlier analysis, the probability that *none* of these safe collapses touches the minimum cut is at least

$$\prod_{i=n/\sqrt{2}+2}^n \frac{i-2}{i} = \frac{(n/\sqrt{2})(n/\sqrt{2}+1)}{n(n-1)} = \frac{n+\sqrt{2}}{2(n-1)} > \frac{1}{2}.$$

Now, to get around the danger of the dangerous phase, we use amplification. However, instead of running through the dangerous phase once, we run it *twice* and keep the best of the two answers. Naturally, we treat the dangerous phase recursively, so we actually obtain a binary recursion tree, which expands as we get closer to the base case, instead of a single path. More formally, the algorithm looks like this:

```
CONTRACT( $G, m$ ):
  for  $i \leftarrow n$  downto  $m$ 
    pick a random edge  $e$  in  $G$ 
     $G \leftarrow G/e$ 
  return  $G$ 
```

```
BETTERGUESS( $G$ ):
  if  $G$  has more than 8 vertices
     $X_1 \leftarrow \text{BETTERGUESS}(\text{CONTRACT}(G, n/\sqrt{2} + 1))$ 
     $X_2 \leftarrow \text{BETTERGUESS}(\text{CONTRACT}(G, n/\sqrt{2} + 1))$ 
    return  $\min\{X_1, X_2\}$ 
  else
    use brute force
```

This might look like we’re just doing the same thing twice, but remember that CONTRACT (and thus BETTERGUESS) is randomized. Each call to CONTRACT contracts an independent random set of edges;  $X_1$  and  $X_2$  are almost always different cuts.

BETTERGUESS correctly returns the minimum cut unless *both* recursive calls return the wrong result.  $X_1$  is the minimum cut of  $G$  if and only if (1) none of the edges of the minimum cut are CONTRACTED and (2) the recursive call to BETTERGUESS returns the minimum cut of the CONTRACTED graph. Thus, if  $P(n)$  denotes the probability that BETTERGUESS returns a minimum cut of an  $n$ -node graph, then  $X_1$  is the minimum cut with probability at least  $1/2 \cdot P(n/\sqrt{2} + 1)$ . The same argument implies that  $X_2$  is the minimum cut with probability at least  $1/2 \cdot P(n/\sqrt{2} + 1)$ . Because these two events are independent, we have the following recurrence, with base case  $P(n) = 1$  for all  $n \leq 6$ .

$$P(n) \geq 1 - \left(1 - \frac{1}{2} P\left(\frac{n}{\sqrt{2}} + 1\right)\right)^2$$

Using a series of transformations, Karger and Stein prove that  $P(n) = \Omega(1/\log n)$ . I’ve included the proof at the end of this note.

For the running time, we get a simple recurrence that is easily solved using recursion trees or the Master theorem (after a domain transformation to remove the  $+1$  from the recurrence).

$$T(n) = O(n^2) + 2T\left(\frac{n}{\sqrt{2}} + 1\right) = O(n^2 \log n)$$

<sup>3</sup>David R. Karger\* and Cliff Stein. An  $\tilde{O}(n^2)$  algorithm for minimum cuts. Proc. 25th STOC, 757–765, 1993.

So all this splitting and recursing has slowed down the guessing algorithm slightly, but the probability of failure is *exponentially* smaller!

Let's express the lower bound  $P(n) = \Omega(1/\log n)$  explicitly as  $P(n) \geq \alpha/\ln n$  for some constant  $\alpha$ . (Karger and Klein's proof implies  $\alpha > 2$ ). If we call BETTERGUESS  $N = c \ln^2 n$  times, for some new constant  $c$ , the overall probability of success is at least

$$1 - \left(1 - \frac{\alpha}{\ln n}\right)^{c \ln^2 n} \geq 1 - e^{-(c/\alpha) \ln n} = 1 - \frac{1}{n^{c/\alpha}}.$$

By setting  $c$  sufficiently large, we can bound the probability of failure by an arbitrarily small polynomial function of  $n$ . In other words, we now have an algorithm that computes the minimum cut with high probability in only  $O(n^2 \log^3 n)$  time!

### \*13.5 Solving the Karger/Stein recurrence

Recall the following recurrence for the probability that BETTERGUESS successfully finds a minimum cut of an  $n$ -node graph:

$$P(n) \geq 1 - \left(1 - \frac{1}{2} P\left(\frac{n}{\sqrt{2}} + 1\right)\right)^2$$

Karger and Stein solve this rather ugly recurrence through a series of functional transformations. Let  $p(k)$  denote the probability of success at the  $k$ th level of recursion, counting upward from the base case. This function satisfies the recurrence

$$p(k) \geq 1 - \left(1 - \frac{p(k-1)}{2}\right)^2 = p(k-1) - \frac{p(k-1)^2}{4}$$

with base case  $p(0) = 1$ . Let  $\bar{p}(k)$  be the function that satisfies this recurrence with equality; clearly,  $p(k) \geq \bar{p}(k)$ . Substituting the function  $z(k) = 4/\bar{p}(k) - 1$  into this recurrence implies (after a bit of algebra) gives a new recurrence

$$z(k) = z(k-1) + 2 + \frac{1}{z(k-1)}$$

with base case  $z(0) = 3$ . Clearly  $z(k) > 1$  for all  $k$ , so we have a conservative upper bound  $z(k) < z(k-1) + 3$ , which implies (by induction) that  $z(k) \leq 3k + 3$ . Substituting  $\bar{p}(k) = 4/(z(k) + 1)$  into this solution, we conclude that

$$p(k) \geq \bar{p}(k) > \frac{1}{3k+6} = \Omega(1/k).$$

To compute the number of levels of recursion that BETTERGUESS executes for an  $n$ -node graph, we solve the secondary recurrence

$$k(n) = 1 + k\left(\frac{n}{\sqrt{2}} + 1\right)$$

with base cases  $k(n) = 0$  for all  $n \leq 8$ . After a domain transformation to remove the  $+1$  from the right side, the recursion tree method (or the Master theorem) implies that  $k(n) = \Theta(\log n)$ .

We conclude that  $P(n) = p(k(n)) = \Omega(1/\log n)$ , as promised.

## Exercises

1. Suppose you had an algorithm to compute the minimum spanning tree of a graph in  $O(m)$  time, where  $m$  is the number of edges in the input graph.<sup>4</sup> Use this algorithm as a subroutine to improve the running time of GUESSMINCUT from  $O(n^2)$  to  $O(m)$ .
2. Suppose you are given a graph  $G$  with weighted edges, and your goal is to find a cut whose total weight (not just number of edges) is smallest.
  - (a) Describe an algorithm to select a random edge of  $G$ , where the probability of choosing edge  $e$  is proportional to the weight of  $e$ .
  - (b) Prove that if you use the algorithm from part (a), instead of choosing edges uniformly at random, the probability that GUESSMINCUT returns a minimum-weight cut is still  $\Omega(1/n^2)$ .
  - (c) What is the running time of your modified GUESSMINCUT algorithm?
3. Prove that GUESSMINCUT returns the *second* smallest cut in its input graph with probability  $\Omega(1/n^3)$ . (The second smallest cut could be significantly larger than the minimum cut.)
4. Consider the following generalization of the BETTERGUESS algorithm, where we pass in a real parameter  $0 < \alpha < 1$  in addition to the graph  $G$ .

```

BETTERGUESS( $G, \alpha$ ):
   $n \leftarrow$  number of vertices in  $G$ 
  if  $n > 8$ 
     $X_1 \leftarrow$  BETTERGUESS( $\text{CONTRACT}(G, n/\alpha), \alpha$ )
     $X_2 \leftarrow$  BETTERGUESS( $\text{CONTRACT}(G, n/\alpha), \alpha$ )
    return  $\min\{X_1, X_2\}$ 
  else
    use brute force
  
```

Assume for this question that the input graph  $G$  has a unique minimum cut.

- (a) What is the running time of the modified algorithm, as a function of  $n$  and  $\alpha$ ? [Hint: Consider the cases  $\alpha < \sqrt{2}$ ,  $\alpha = \sqrt{2}$ , and  $\alpha > \sqrt{2}$  separately.]
- (b) What is the probability that  $\text{CONTRACT}(G, n/\alpha)$  does not contract any edge in the minimum cut in  $G$ ? Give both an exact expression involving both  $n$  and  $\alpha$ , and a simple approximation in terms of just  $\alpha$ . [Hint: When  $\alpha = \sqrt{2}$ , the probability is approximately 1/2.]
- (c) Estimate the probability that BETTERGUESS( $G, \alpha$ ) returns the minimum cut in  $G$ , by adapting the solution to the Karger-Stein recurrence. [Hint: Consider the cases  $\alpha < \sqrt{2}$ ,  $\alpha = \sqrt{2}$ , and  $\alpha > \sqrt{2}$  separately.]
- (d) Suppose we iterate BETTERGUESS( $G, \alpha$ ) until we are guaranteed to see the minimum cut with high probability. What is the running time of the resulting algorithm? For which value of  $\alpha$  is this running time minimized?
- (e) Suppose we modify BETTERGUESS( $G, \alpha$ ) further, to recurse four times instead of only twice. Now what is the best choice of  $\alpha$ ? What is the resulting running time?

---

<sup>4</sup>In fact, there is a randomized MST algorithm (due to Philip Klein, David Karger, and Robert Tarjan) whose expected running time is  $O(m)$ .

*The goode workes that men don whil they ben in good lif al amortised by synne folwyng.*

— Geoffrey Chaucer, “The Persones [Parson’s] Tale” (c.1400)

*I will gladly pay you Tuesday for a hamburger today.*

— J. Wellington Wimpy, “Thimble Theatre” (1931)

*I want my two dollars!*

— Johnny Gasparini [Demian Slade], “Better Off Dead” (1985)

## 14 Amortized Analysis

### 14.1 Incrementing a Binary Counter

It is a straightforward exercise in induction, which often appears on Homework 0, to prove that any non-negative integer  $n$  can be represented as the sum of distinct powers of 2. Although some students correctly use induction on the number of bits—pulling off either the least significant bit or the most significant bit in the binary representation and letting the Recursion Fairy convert the remainder—the most commonly submitted proof uses induction on the value of the integer, as follows:

**Proof:** The base case  $n = 0$  is trivial. For any  $n > 0$ , the inductive hypothesis implies that there is set of distinct powers of 2 whose sum is  $n - 1$ . If we add  $2^0$  to this set, we obtain a *multiset* of powers of two whose sum is  $n$ , which might contain two copies of  $2^0$ . Then as long as there are two copies of any  $2^i$  in the multiset, we remove them both and insert  $2^{i+1}$  in their place. The sum of the elements of the multiset is unchanged by this replacement, because  $2^{i+1} = 2^i + 2^i$ . Each iteration decreases the size of the multiset by 1, so the replacement process must eventually terminate. When it does terminate, we have a *set* of distinct powers of 2 whose sum is  $n$ .  $\square$

This proof is describing an algorithm to increment a binary counter from  $n - 1$  to  $n$ . Here’s a more formal (and shorter!) description of the algorithm to add 1 to a binary counter. The input  $B$  is an (infinite) array of bits, where  $B[i] = 1$  if and only if  $2^i$  appears in the sum.

```

INCREMENT( $B[0..\infty]$ ):
   $i \leftarrow 0$ 
  while  $B[i] = 1$ 
     $B[i] \leftarrow 0$ 
     $i \leftarrow i + 1$ 
   $B[i] \leftarrow 1$ 

```

We’ve already argued that INCREMENT must terminate, but how quickly? Obviously, the running time depends on the array of bits passed as input. If the first  $k$  bits are all 1s, then INCREMENT takes  $\Theta(k)$  time. The binary representation of any positive integer  $n$  is exactly  $\lfloor \lg n \rfloor + 1$  bits long. Thus, if  $B$  represents an integer between 0 and  $n$ , INCREMENT takes  $\Theta(\log n)$  time in the worst case.

### 14.2 Counting from 0 to $n$

Now suppose we want to use INCREMENT to count from 0 to  $n$ . If we only use the worst-case running time for each INCREMENT, we get an upper bound of  $O(n \log n)$  on the total running time. Although this bound is correct, we can do better. There are several general methods for proving that the total running time is actually only  $O(n)$ . Many of these methods are logically equivalent, but different formulations are more natural for different problems.

### 14.2.1 The Summation (Aggregate) Method

The easiest way to get a tighter bound is to observe that we don't need to flip  $\Theta(\log n)$  bits every time INCREMENT is called. The least significant bit  $B[0]$  does flip every time, but  $B[1]$  only flips every other time,  $B[2]$  flips every 4th time, and in general,  $B[i]$  flips every  $2^i$ th time. If we start with an array full of 0's, a sequence of  $n$  INCREMENTS flips each bit  $B[i]$  exactly  $\lfloor n/2^i \rfloor$  times. Thus, the total number of bit-flips for the entire sequence is

$$\sum_{i=0}^{\lfloor \lg n \rfloor} \left\lfloor \frac{n}{2^i} \right\rfloor < \sum_{i=0}^{\infty} \frac{n}{2^i} = 2n.$$

Thus, *on average*, each call to INCREMENT flips only two bits, and so runs in constant time.

This 'on average' is quite different from the averaging we consider with randomized algorithms. There is no probability involved; we are averaging over a sequence of operations, not the possible running times of a single operation. This averaging idea is called *amortization*—the *amortized* time for each INCREMENT is  $O(1)$ . Amortization is a *sleazy* clever trick used by accountants to average large one-time costs over long periods of time; the most common example is calculating uniform payments for a loan, even though the borrower is paying interest on less and less capital over time. For this reason, 'amortized cost' is a common synonym for amortized running time.

Most textbooks call this technique for bounding amortized costs the *aggregate* method, or *aggregate analysis*, but this is just a fancy name for adding up the costs of the individual operations and dividing by the number of operations.

**The Summation Method.** *Let  $T(n)$  be the worst-case running time for a sequence of  $n$  operations. The amortized time for each operation is  $T(n)/n$ .*

### 14.2.2 The Taxation (Accounting) Method

A second method we can use to derive amortized bounds is called either the *accounting* method or the *taxation* method. Suppose it costs us a dollar to toggle a bit, so we can measure the running time of our algorithm in dollars. Time is money!

Instead of paying for each bit flip when it happens, the Increment Revenue Service charges a two-dollar *increment tax* whenever we want to set a bit from zero to one. One of those dollars is spent changing the bit from zero to one; the other is stored away as *credit* until we need to reset the same bit to zero. The key point here is that we always have enough credit saved up to pay for the next INCREMENT. The amortized cost of an INCREMENT is the total tax it incurs, which is exactly 2 dollars, since each INCREMENT changes just one bit from 0 to 1.

It is often useful to distribute the tax income to specific pieces of the data structure. For example, for each INCREMENT, we could store one of the two dollars on the single bit that is set for 0 to 1, so that *that* bit can pay to reset itself back to zero later on.

**Taxation Method 1.** *Certain steps in the algorithm charge you taxes, so that the total cost incurred by the algorithm is never more than the total tax you pay. The amortized cost of an operation is the overall tax charged to you during that operation.*

A different way to schedule the taxes is for *every* bit to charge us a tax at *every* operation, regardless of whether the bit changes or not. Specifically, each bit  $B[i]$  charges a tax of  $1/2^i$  dollars for each INCREMENT. The total tax we are charged during each INCREMENT is  $\sum_{i \geq 0} 2^{-i} = 2$  dollars. Every time a bit  $B[i]$  actually needs to be flipped, it has collected exactly \$1, which is just enough for us to pay for the flip.

**Taxation Method 2.** *Certain portions of the data structure charge you taxes at each operation, so that the total cost of maintaining the data structure is never more than the total taxes you pay. The amortized cost of an operation is the overall tax you pay during that operation.*

In both of the taxation methods, our task as algorithm analysts is to come up with an appropriate ‘tax schedule’. Different ‘schedules’ can result in different amortized time bounds. The tightest bounds are obtained from tax schedules that just barely stay in the black.

#### 14.2.3 The Charging Method

Another common method of amortized analysis involves *charging* the cost of some steps to some other, earlier steps. The method is similar to taxation, except that we focus on where each unit of tax is (or will be) spent, rather than where it is collected. By charging the cost of some operations to earlier operations, we are overestimating the total cost of any sequence of operations, since we pay for some charges from future operations that may never actually occur.

For example, in our binary counter, suppose we charge the cost of clearing a bit (changing its value from 1 to 0) to the previous operation that sets that bit (changing its value from 0 to 1). If we flip  $k$  bits during an INCREMENT, we charge  $k - 1$  of those bit-flips to earlier bit-flips. Conversely, the single operation that sets a bit receives at most one unit of charge from the next time that bit is cleared. So instead of paying for  $k$  bit-flips, we pay for at most two: one for actually setting a bit, plus at most one charge from the future for clearing that same bit. Thus, the total amortized cost of the INCREMENT is at most two bit-flips.

**The Charging Method.** *Charge the cost of some steps of the algorithm to earlier steps, or to steps in some earlier operation. The amortized cost of the algorithm is its actual running time, minus its total charges to past operations, plus its total charge from future operations.*

#### 14.2.4 The Potential Method

The most powerful method (and the hardest to use) builds on a physics metaphor of ‘potential energy’. Instead of associating costs or taxes with particular operations or pieces of the data structure, we represent prepaid work as *potential* that can be spent on later operations. The potential is a function of the entire data structure.

Let  $D_i$  denote our data structure after  $i$  operations, and let  $\Phi_i$  denote its potential. Let  $c_i$  denote the actual cost of the  $i$ th operation (which changes  $D_{i-1}$  into  $D_i$ ). Then the *amortized* cost of the  $i$ th operation, denoted  $a_i$ , is defined to be the actual cost plus the increase in potential:

$$a_i = c_i + \Phi_i - \Phi_{i-1}$$

So the *total* amortized cost of  $n$  operations is the actual total cost plus the total increase in potential:

$$\sum_{i=1}^n a_i = \sum_{i=1}^n (c_i + \Phi_i - \Phi_{i-1}) = \sum_{i=1}^n c_i + \Phi_n - \Phi_0.$$

A potential function is *valid* if  $\Phi_i - \Phi_0 \geq 0$  for all  $i$ . If the potential function is valid, then the total *actual* cost of any sequence of operations is always less than the total *amortized* cost:

$$\sum_{i=1}^n c_i = \sum_{i=1}^n a_i - \Phi_n \leq \sum_{i=1}^n a_i.$$

For our binary counter example, we can define the potential  $\Phi_i$  after the  $i$ th INCREMENT to be the number of bits with value 1. Initially, all bits are equal to zero, so  $\Phi_0 = 0$ , and clearly  $\Phi_i > 0$  for all  $i > 0$ , so this is a valid potential function. We can describe both the actual cost of an INCREMENT and the change in potential in terms of the number of bits set to 1 and reset to 0.

$$c_i = \#\text{bits changed from 0 to 1} + \#\text{bits changed from 1 to 0}$$

$$\Phi_i - \Phi_{i-1} = \#\text{bits changed from 0 to 1} - \#\text{bits changed from 1 to 0}$$

Thus, the amortized cost of the  $i$ th INCREMENT is

$$a_i = c_i + \Phi_i - \Phi_{i-1} = 2 \times \#\text{bits changed from 0 to 1}$$

Since INCREMENT changes only *one* bit from 0 to 1, the amortized cost INCREMENT is 2.

**The Potential Method.** Define a potential function for the data structure that is initially equal to zero and is always nonnegative. The amortized cost of an operation is its actual cost plus the change in potential.

For this particular example, the potential is precisely the total unspent taxes paid using the taxation method, so not too surprisingly, we obtain precisely the same amortized cost. In general, however, there may be no way of interpreting the change in potential as ‘taxes’. Taxation and charging are useful when there is a convenient way to distribute costs to specific steps in the algorithm or components of the data structure; potential arguments allow us to argue more globally when a local distribution is difficult or impossible.

Different potential functions can lead to different amortized time bounds. The trick to using the potential method is to come up with the best possible potential function. A good potential function goes up a little during any cheap/fast operation, and goes down a lot during any expensive/slow operation. Unfortunately, there is no general technique for doing this other than playing around with the data structure and trying lots of different possibilities.

### 14.3 Incrementing and Decrementing

Now suppose we wanted a binary counter that we can both increment and decrement efficiently. A standard binary counter won’t work, even in an amortized sense; if we alternate between  $2^k$  and  $2^k - 1$ , every operation costs  $\Theta(k)$  time.

A nice alternative is represent a number as a pair of bit strings  $(P, N)$ , where for any bit position  $i$ , at most one of the bits  $P[i]$  and  $N[i]$  is equal to 1. The actual value of the counter is  $P - N$ . Here are algorithms to increment and decrement our double binary counter.

<u>INCREMENT(<math>P, N</math>):</u>	<u>DECREMENT(<math>P, N</math>):</u>
$i \leftarrow 0$ while $P[i] = 1$ $P[i] \leftarrow 0$ $i \leftarrow i + 1$ if $N[i] = 1$ $N[i] \leftarrow 0$ else $P[i] \leftarrow 1$	$i \leftarrow 0$ while $N[i] = 1$ $N[i] \leftarrow 0$ $i \leftarrow i + 1$ if $P[i] = 1$ $P[i] \leftarrow 0$ else $N[i] \leftarrow 1$

Here’s an example of these algorithms in action. Notice that any number other than zero can be represented in multiple (in fact, infinitely many) ways.

$$\begin{array}{ccccccc}
 P = 10001 & P = 100\textcolor{red}{10} & P = 10011 & P = 100\textcolor{red}{00} & P = 10000 & P = 10000 & P = 10001 \\
 N = 01100 & \xrightarrow{++} N = 01100 & \xrightarrow{++} N = 01100 & \xrightarrow{++} N = 01000 & \xrightarrow{--} N = 01001 & \xrightarrow{--} N = 01010 & \xrightarrow{++} N = 01010 \\
 P - N = 5 & P - N = 6 & P - N = 7 & P - N = 8 & P - N = 7 & P - N = 6 & P - N = 7
 \end{array}$$

Incrementing and decrementing a double-binary counter.

Now suppose we start from  $(0, 0)$  and apply a sequence of  $n$  INCREMENTS and DECREMENTS. In the worst case, operation takes  $\Theta(\log n)$  time, but what is the amortized cost? We can't use the aggregate method here, since we don't know what the sequence of operations looks like.

What about the taxation method? It's not hard to prove (by induction, of course) that after either  $P[i]$  or  $N[i]$  is set to 1, there must be at least  $2^i$  operations, either INCREMENTS or DECREMENTS, before that bit is reset to 0. So if each bit  $P[i]$  and  $N[i]$  pays a tax of  $2^{-i}$  at each operation, we will always have enough money to pay for the next operation. Thus, the amortized cost of each operation is at most  $\sum_{i \geq 0} 2 \cdot 2^{-i} = 4$ .

We can get even better bounds using the potential method. Define the potential  $\Phi_i$  to be the number of 1-bits in both  $P$  and  $N$  after  $i$  operations. Just as before, we have

$$\begin{aligned}
 c_i &= \# \text{bits changed from 0 to 1} + \# \text{bits changed from 1 to 0} \\
 \Phi_i - \Phi_{i-1} &= \# \text{bits changed from 0 to 1} - \# \text{bits changed from 1 to 0} \\
 \implies a_i &= 2 \times \# \text{bits changed from 0 to 1}
 \end{aligned}$$

Since each operation changes *at most* one bit to 1, the  $i$ th operation has amortized cost  $a_i \leq 2$ .

#### \*14.4 Gray Codes

An attractive alternate solution to the increment/decrement problem was independently suggested by several students. *Gray codes* (named after Frank Gray, who discovered them in the 1950s) are methods for representing numbers as bit strings so that successive numbers differ by only one bit. For example, here is the four-bit *binary reflected* Gray code for the integers 0 through 15:

0000, 0001, 0011, 0010, 0110, 0111, 0101, 0100, 1100, 1101, 1111, 1110, 1010, 1011, 1001, 1000

The general rule for incrementing a binary reflected Gray code is to invert the bit that would be set from 0 to 1 by a normal binary counter. In terms of bit-flips, this is the perfect solution; each increment or decrement *by definition* changes only one bit. Unfortunately, the naïve algorithm to *find* the single bit to flip still requires  $\Theta(\log n)$  time in the worst case. Thus, so the total cost of maintaining a Gray code, using the obvious algorithm, is the same as that of maintaining a normal binary counter.

Fortunately, this is only true of the naïve algorithm. The following algorithm, discovered by Gideon Ehrlich<sup>1</sup> in 1973, maintains a Gray code counter in constant *worst-case* time per increment! The algorithm uses a separate ‘focus’ array  $F[0..n]$  in addition to a Gray-code bit array  $G[0..n-1]$ .

```

EHRLICHGRAYINIT( $n$ ):
  for  $i \leftarrow 0$  to  $n - 1$ 
     $G[i] \leftarrow 0$ 
  for  $i \leftarrow 0$  to  $n$ 
     $F[i] \leftarrow i$ 

```

```

EHRLICHGRAYINCREMENT( $n$ ):
   $j \leftarrow F[0]$ 
   $F[0] \leftarrow 0$ 
  if  $j = n$ 
     $G[n-1] \leftarrow 1 - G[n-1]$ 
  else
     $G[j] = 1 - G[j]$ 
     $F[j] \leftarrow F[j+1]$ 
     $F[j+1] \leftarrow j+1$ 

```

<sup>1</sup>Gideon Ehrlich. Loopless algorithms for generating permutations, combinations, and other combinatorial configurations. *J. Assoc. Comput. Mach.* 20:500–513, 1973.

The EHRLICHGRAYINCREMENT algorithm obviously runs in  $O(1)$  time, even in the worst case. Here's the algorithm in action with  $n = 4$ . The first line is the Gray bit-vector  $G$ , and the second line shows the focus vector  $F$ , both in reverse order:

```
G : 0000, 0001, 0011, 0010, 0110, 0111, 0101, 0100, 1100, 1101, 1111, 1110, 1010, 1011, 1001, 1000
F : 3210, 3211, 3220, 3212, 3310, 3311, 3230, 3213, 4210, 4211, 4220, 4212, 3410, 3411, 3240, 3214
```

Voodoo! I won't explain in detail how Ehrlich's algorithm works, except to point out the following invariant. Let  $B[i]$  denote the  $i$ th bit in the *standard* binary representation of the current number. **If  $B[j] = 0$  and  $B[j - 1] = 1$ , then  $F[j]$  is the smallest integer  $k > j$  such that  $B[k] = 1$ ; otherwise,  $F[j] = j$ .** Got that?

But wait — this algorithm only handles increments; what if we also want to decrement? Sorry, I don't have a clue. Extra credit, anyone?

## 14.5 Generalities and Warnings

Although computer scientists usually apply amortized analysis to understand the efficiency of maintaining and querying data structures, you should remember that amortization can be applied to *any* sequence of numbers. Banks have been using amortization to calculate fixed payments for interest-bearing loans for centuries. The IRS allows taxpayers to amortize business expenses or gambling losses across several years for purposes of computing income taxes. Some cell phone contracts let you to apply amortization to calling time, by rolling unused minutes from one month into the next month.

It's also important to remember that ***amortized time bounds are not unique***. For a data structure that supports multiple operations, different amortization schemes can assign different costs to *exactly the same* algorithms. For example, consider a generic data structure that can be modified by three algorithms: FOLD, SPINDLE, and MUTILATE. One amortization scheme might imply that FOLD and SPINDLE each run in  $O(\log n)$  amortized time, while MUTILATE runs in  $O(n)$  amortized time. Another scheme might imply that FOLD runs in  $O(\sqrt{n})$  amortized time, while SPINDLE and MUTILATE each run in  $O(1)$  amortized time. These two results are not necessarily inconsistent! Moreover, there is no general reason to prefer one of these sets of amortized time bounds over the other; our preference may depend on the context in which the data structure is used.

## Exercises

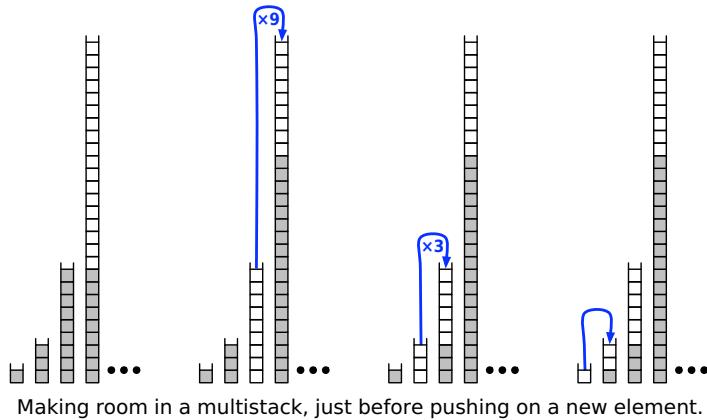
1. Suppose we are maintaining a data structure under a series of operations. Let  $f(n)$  denote the actual running time of the  $n$ th operation. For each of the following functions  $f$ , determine the resulting amortized cost of a single operation. (For practice, try *all* of the methods described in this note.)
  - (a)  $f(n)$  is the largest integer  $i$  such that  $2^i$  divides  $n$ .
  - (b)  $f(n)$  is the largest power of 2 that divides  $n$ .
  - (c)  $f(n) = n$  if  $n$  is a power of 2, and  $f(n) = 1$  otherwise.
  - (d)  $f(n) = n^2$  if  $n$  is a power of 2, and  $f(n) = 1$  otherwise.
2. A *multistack* consists of an infinite series of stacks  $S_0, S_1, S_2, \dots$ , where the  $i$ th stack  $S_i$  can hold up to  $3^i$  elements. The user always pushes and pops elements from the smallest stack  $S_0$ . However, before any element can be pushed onto any full stack  $S_i$ , we first pop all the elements off  $S_i$  and

push them onto stack  $S_{i+1}$  to make room. (Thus, if  $S_{i+1}$  is already full, we first recursively move all its members to  $S_{i+2}$ .) Similarly, before any element can be popped from any empty stack  $S_i$ , we first pop  $3^i$  elements from  $S_{i+1}$  and push them onto  $S_i$  to make room. (Thus, if  $S_{i+1}$  is already empty, we first recursively fill it by popping elements from  $S_{i+2}$ .) Moving a single element from one stack to another takes  $O(1)$  time.

Here is pseudocode for the multistack operations `MSPUSH` and `MSPOP`. The internal stacks are managed with the subroutines `PUSH` and `POP`.

```
MPUSH( $x$ ):
 $i \leftarrow 0$ 
while  $S_i$  is full
     $i \leftarrow i + 1$ 
while  $i > 0$ 
     $i \leftarrow i - 1$ 
    for  $j \leftarrow 1$  to  $3^i$ 
        PUSH( $S_{i+1}$ , POP( $S_i$ ))
PUSH( $S_0, x$ )
```

```
MPOP( $x$ ):
 $i \leftarrow 0$ 
while  $S_i$  is empty
     $i \leftarrow i + 1$ 
while  $i > 0$ 
     $i \leftarrow i - 1$ 
    for  $j \leftarrow 1$  to  $3^i$ 
        PUSH( $S_i$ , POP( $S_{i+1}$ ))
return POP( $S_0$ )
```



- (a) In the worst case, how long does it take to push one more element onto a multistack containing  $n$  elements?
- (b) Prove that if the user never pops anything from the multistack, the amortized cost of a push operation is  $O(\log n)$ , where  $n$  is the maximum number of elements in the multistack during its lifetime.
- (c) Prove that in any intermixed sequence of pushes and pops, each push or pop operation takes  $O(\log n)$  amortized time, where  $n$  is the maximum number of elements in the multistack during its lifetime.

3. Remember the difference between stacks and queues? Good.

- (a) Describe how to implement a queue using two stacks and  $O(1)$  additional memory, so that the amortized time for any enqueue or dequeue operation is  $O(1)$ . The *only* access you have to the stacks is through the standard subroutines `PUSH` and `POP`.
- (b) A *quack* is a data structure combining properties of both stacks and queues. It can be viewed as a list of elements written left to right such that three operations are possible:

- **Push:** add a new item to the left end of the list;
- **Pop:** remove the item on the left end of the list;
- **Pull:** remove the item on the right end of the list.

Implement a quack using *three* stacks and  $O(1)$  additional memory, so that the amortized time for any push, pop, or pull operation is  $O(1)$ . In particular, each element pushed onto the quack should be stored in exactly one of the three stacks. Again, you are *only* allowed to access the stacks through the standard functions **PUSH** and **POP**.

4. Suppose we can insert or delete an element into a hash table in  $O(1)$  time. In order to ensure that our hash table is always big enough, without wasting a lot of memory, we will use the following global rebuilding rules:

- After an insertion, if the table is more than  $3/4$  full, we allocate a new table twice as big as our current table, insert everything into the new table, and then free the old table.
- After a deletion, if the table is less than  $1/4$  full, we allocate a new table half as big as our current table, insert everything into the new table, and then free the old table.

Show that for any sequence of insertions and deletions, the amortized time per operation is still  $O(1)$ . [*Hint: Do not even look at the potential-function argument in CLRS; there is a much easier solution!*]

5. Chicago has many tall buildings, but only some of them have a clear view of Lake Michigan. Suppose we are given an array  $A[1..n]$  that stores the height of  $n$  buildings on a city block, indexed from west to east. Building  $i$  has a good view of Lake Michigan if and only if every building to the east of  $i$  is shorter than  $i$ .

Here is an algorithm that computes which buildings have a good view of Lake Michigan. What is the running time of this algorithm?

```

GOODVIEW( $A[1..n]$ ):
  initialize a stack  $S$ 
  for  $i \leftarrow 1$  to  $n$ 
    while ( $S$  not empty and  $A[i] > A[\text{TOP}(S)]$ )
      POP( $S$ )
      PUSH( $S, i$ )
  return  $S$ 
```

6. Design and analyze a simple data structure that maintains a list of integers and supports the following operations.

- **CREATE()** creates and returns a new list
- **PUSH( $L, x$ )** appends  $x$  to the end of  $L$
- **POP( $L$ )** deletes the last entry of  $L$  and returns it
- **LOOKUP( $L, k$ )** returns the  $k$ th entry of  $L$

Your solution may use these primitive data structures: arrays, balanced binary search trees, heaps, queues, single or doubly linked lists, and stacks. If your algorithm uses *anything* fancier, you must give an explicit implementation. Your data structure must support all operations in amortized

constant time. In addition, your data structure must support each LOOKUP in *worst-case*  $O(1)$  time. At all times, the size of your data structure must be linear in the number of objects it stores.

- \*7. Suppose instead of powers of two, we represent integers as the sum of Fibonacci numbers. In other words, instead of an array of bits, we keep an array of *fits*, where the  $i$ th least significant fit indicates whether the sum includes the  $i$ th Fibonacci number  $F_i$ . For example, the fitstring  $101110_F$  represents the number  $F_6 + F_4 + F_3 + F_2 = 8 + 3 + 2 + 1 = 14$ . Describe algorithms to increment and decrement a single fitstring in constant amortized time. [*Hint: Most numbers can be represented by more than one fitstring!*]
- \*8. A *doubly lazy binary counter* represents any number as a weighted sum of powers of two, where each weight is one of four values:  $-1$ ,  $0$ ,  $1$ , or  $2$ . (For succinctness, I'll write  $\pm$  instead of  $-1$ .) Every integer—positive, negative, or zero—has an infinite number of doubly lazy binary representations. For example, the number  $13$  can be represented as  $1101$  (the standard binary representation), or  $2\pm 01$  (because  $2 \cdot 2^3 - 2^2 + 2^0 = 13$ ) or  $10\pm 1\pm$  (because  $2^4 - 2^2 + 2^1 - 2^0 = 13$ ) or  $\pm 1200010\pm 1\pm$  (because  $-2^{10} + 2^9 + 2 \cdot 2^8 + 2^4 - 2^2 + 2^1 - 2^0 = 13$ ).

To increment a doubly lazy binary counter, we add 1 to the least significant bit, then carry the rightmost 2 (if any). To decrement, we subtract 1 from the least significant bit, and then borrow the rightmost  $\pm$  (if any).

<u>LAZYINCREMENT(<math>B[0..n]</math>):</u>
$B[0] \leftarrow B[0] + 1$
for $i \leftarrow 1$ to $n - 1$
if $B[i] = 2$
$B[i] \leftarrow 0$
$B[i + 1] \leftarrow B[i + 1] + 1$
return

<u>LAZYDECREMENT(<math>B[0..n]</math>):</u>
$B[0] \leftarrow B[0] - 1$
for $i \leftarrow 1$ to $n - 1$
if $B[i] = -1$
$B[i] \leftarrow 1$
$B[i + 1] \leftarrow B[i + 1] - 1$
return

For example, here is a doubly lazy binary count from zero up to twenty and then back down to zero. The bits are written with the least significant bit  $B[0]$  on the *right*, omitting all leading 0's

$0 \xrightarrow{++} 1 \xrightarrow{++} 10 \xrightarrow{++} 11 \xrightarrow{++} 20 \xrightarrow{++} 101 \xrightarrow{++} 110 \xrightarrow{++} 111 \xrightarrow{++} 120 \xrightarrow{++} 201 \xrightarrow{++} 210$   
 $\xrightarrow{++} 1011 \xrightarrow{++} 1020 \xrightarrow{++} 1101 \xrightarrow{++} 1110 \xrightarrow{++} 1111 \xrightarrow{++} 1120 \xrightarrow{++} 1201 \xrightarrow{++} 1210 \xrightarrow{++} 2011 \xrightarrow{++} 2020$   
 $\xrightarrow{--} 2011 \xrightarrow{--} 2010 \xrightarrow{--} 2001 \xrightarrow{--} 2000 \xrightarrow{--} 20\pm 1 \xrightarrow{--} 2\pm 10 \xrightarrow{--} 2\pm 01 \xrightarrow{--} 1100 \xrightarrow{--} 11\pm 1 \xrightarrow{--} 1010$   
 $\xrightarrow{--} 1001 \xrightarrow{--} 1000 \xrightarrow{--} 10\pm 1 \xrightarrow{--} 1\pm 10 \xrightarrow{--} 1\pm 01 \xrightarrow{--} 100 \xrightarrow{--} 1\pm 1 \xrightarrow{--} 10 \xrightarrow{--} 1 \xrightarrow{--} 0$

Prove that for any intermixed sequence of increments and decrements of a doubly lazy binary number, starting with 0, the amortized time for each operation is  $O(1)$ . Do *not* assume, as in the example above, that all the increments come before all the decrements.

*Everything was balanced before the computers went off line. Try and adjust something, and you unbalance something else. Try and adjust that, you unbalance two more and before you know what's happened, the ship is out of control.*

— Blake, *Blake's 7, "Breakdown"* (March 6, 1978)

*A good scapegoat is nearly as welcome as a solution to the problem.*

— Anonymous

*Let's play.*

— El Mariachi [Antonio Banderas], *Desperado* (1992)

*CAPTAIN: TAKE OFF EVERY 'ZIG'!!*

*CAPTAIN: YOU KNOW WHAT YOU DOING.*

*CAPTAIN: MOVE 'ZIG'.*

*CAPTAIN: FOR GREAT JUSTICE.*

— *Zero Wing* (1992)

## 15 Scapegoat and Splay Trees

### 15.1 Definitions

I'll assume that everyone is already familiar with the standard terminology for binary search trees—node, search key, edge, root, internal node, leaf, right child, left child, parent, descendant, sibling, ancestor, subtree, preorder, postorder, inorder, etc.—as well as the standard algorithms for searching for a node, inserting a node, or deleting a node. Otherwise, consult your favorite data structures textbook.

For this lecture, we will consider only *full* binary trees—where every internal node has *exactly* two children—where only the *internal* nodes actually store search keys. In practice, we can represent the leaves with null pointers.

Recall that the *depth* of a node is its distance from the root, and its *height* is the distance to the farthest leaf in its subtree. The height (or depth) of the tree is just the height of the root. The *size* of a node is the number of nodes in its subtree. The size  $n$  of the whole tree is just the total number of nodes.

A tree with height  $h$  has at most  $2^h$  leaves, so the minimum height of an  $n$ -leaf binary tree is  $\lceil \lg n \rceil$ . In the worst case, the time required for a search, insertion, or deletion to the height of the tree, so in general we would like keep the height as close to  $\lg n$  as possible. The best we can possibly do is to have a *perfectly balanced* tree, in which each subtree has as close to half the leaves as possible, and both subtrees are perfectly balanced. The height of a perfectly balanced tree is  $\lceil \lg n \rceil$ , so the worst-case search time is  $O(\log n)$ . However, even if we started with a perfectly balanced tree, a malicious sequence of insertions and/or deletions could make the tree arbitrarily unbalanced, driving the search time up to  $\Theta(n)$ .

To avoid this problem, we need to periodically modify the tree to maintain ‘balance’. There are several methods for doing this, and depending on the method we use, the search tree is given a different name. Examples include AVL trees, red-black trees, height-balanced trees, weight-balanced trees, bounded-balance trees, path-balanced trees,  $B$ -trees, treaps, randomized binary search trees, skip lists,<sup>1</sup> and jumplists. Some of these trees support searches, insertions, and deletions, in  $O(\log n)$  *worst-case* time, others in  $O(\log n)$  *amortized* time, still others in  $O(\log n)$  *expected* time.

In this lecture, I'll discuss two binary search tree data structures with good *amortized* performance. The first is the *scapegoat tree*, discovered by Arne Andersson\* in 1989 [1, 2] and independently<sup>2</sup> by Igal

<sup>1</sup>Yeah, yeah. Skip lists aren't really binary search trees. Whatever you say, Mr. Picky.

<sup>2</sup>The claim of independence is Andersson's [2]. The two papers actually describe very slightly different rebalancing algorithms. The algorithm I'm using here is closer to Andersson's, but my analysis is closer to Galperin and Rivest's.

Galperin\* and Ron Rivest in 1993 [10]. The second is the *splay tree*, discovered by Danny Sleator and Bob Tarjan in 1981 [16, 14].

## 15.2 Lazy Deletions: Global Rebuilding

First let's consider the simple case where we start with a perfectly-balanced tree, and we only want to perform searches and deletions. To get good search and delete times, we will use a technique called *global rebuilding*. When we get a delete request, we locate and mark the node to be deleted, *but we don't actually delete it*. This requires a simple modification to our search algorithm—we still use marked nodes to guide searches, but if we search for a marked node, the search routine says it isn't there. This keeps the tree more or less balanced, but now the search time is no longer a function of the amount of data currently stored in the tree. To remedy this, we also keep track of how many nodes have been marked, and then apply the following rule:

**Global Rebuilding Rule.** *As soon as half the nodes in the tree have been marked, rebuild a new perfectly balanced tree containing only the unmarked nodes.*<sup>3</sup>

With this rule in place, a search takes  $O(\log n)$  time in the worst case, where  $n$  is the number of unmarked nodes. Specifically, since the tree has at most  $n$  marked nodes, or  $2n$  nodes altogether, we need to examine at most  $\lg n + 1$  keys. There are several methods for rebuilding the tree in  $O(n)$  time, where  $n$  is the size of the new tree. (Homework!) So a single deletion can cost  $\Theta(n)$  time in the worst case, but only if we have to rebuild; most deletions take only  $O(\log n)$  time.

We spend  $O(n)$  time rebuilding, but only after  $\Omega(n)$  deletions, so the *amortized* cost of rebuilding the tree is  $O(1)$  per deletion. (Here I'm using a simple version of the 'taxation method'. For each deletion, we charge a \$1 tax; after  $n$  deletions, we've collected \$ $n$ , which is just enough to pay for rebalancing the tree containing the remaining  $n$  nodes.) Since we also have to find and mark the node being 'deleted', the total amortized time for a deletion is  $O(\log n)$ .

## 15.3 Insertions: Partial Rebuilding

Now suppose we only want to support searches and insertions. We can't 'not really insert' new nodes into the tree, since that would make them unavailable to the search algorithm.<sup>4</sup> So instead, we'll use another method called *partial rebuilding*. We will insert new nodes normally, but whenever a *subtree* becomes unbalanced enough, we rebuild it. The definition of 'unbalanced enough' depends on an arbitrary constant  $\alpha > 1$ .

Each node  $v$  will now also store  $height(v)$  and  $size(v)$ . We now modify our insertion algorithm with the following rule:

**Partial Rebuilding Rule.** *After we insert a node, walk back up the tree updating the heights and sizes of the nodes on the search path. If we encounter a node  $v$  where  $height(v) > \alpha \cdot \lg(size(v))$ , rebuild its subtree into a perfectly balanced tree (in  $O(size(v))$  time).*

If we always follow this rule, then after an insertion, the height of the tree is at most  $\alpha \cdot \lg n$ . Thus, since  $\alpha$  is a constant, the worst-case search time is  $O(\log n)$ . In the worst case, insertions require  $\Theta(n)$  time—we might have to rebuild the entire tree. However, the *amortized* time for each insertion is again only  $O(\log n)$ . Not surprisingly, the proof is a little bit more complicated than for deletions.

---

<sup>3</sup>Alternately: When the number of unmarked nodes is one less than an exact power of two, rebuild the tree. This rule ensures that the tree is always *exactly* balanced.

<sup>4</sup>Well, we could use the Bentley-Saxe\* logarithmic method [3], but that would raise the query time to  $O(\log^2 n)$ .

Define the *imbalance*  $I(v)$  of a node  $v$  to be one less than the absolute difference between the sizes of its two subtrees, or zero, whichever is larger:

$$I(v) := \max \{0, |\text{size}(\text{left}(v)) - \text{size}(\text{right}(v))| - 1\}$$

A simple induction proof implies that  $I(v) = 0$  for every node  $v$  in a perfectly balanced tree. So immediately after we rebuild the subtree of  $v$ , we have  $I(v) = 0$ . On the other hand, each insertion into the subtree of  $v$  increments either  $\text{size}(\text{left}(v))$  or  $\text{size}(\text{right}(v))$ , so  $I(v)$  changes by at most 1.

The whole analysis boils down to the following lemma.

**Lemma 1.** *Just before we rebuild  $v$ 's subtree,  $I(v) = \Omega(\text{size}(v))$ .*

Before we prove this, let's first look at what it implies. If  $I(v) = \Omega(\text{size}(v))$ , then  $\Omega(\text{size}(v))$  keys have been inserted in the  $v$ 's subtree since the last time it was rebuilt from scratch. On the other hand, rebuilding the subtree requires  $O(\text{size}(v))$  time. Thus, if we amortize the rebuilding cost across all the insertions since the last rebuilding,  $v$  is charged *constant* time for each insertion into its subtree. Since each new key is inserted into at most  $\alpha \cdot \lg n = O(\log n)$  subtrees, the total amortized cost of an insertion is  $O(\log n)$ .

**Proof:** Since we're about to rebuild the subtree at  $v$ , we must have  $\text{height}(v) > \alpha \cdot \lg \text{size}(v)$ . Without loss of generality, suppose that the node we just inserted went into  $v$ 's left subtree. Either we just rebuilt this subtree or we didn't have to, so we also have  $\text{height}(\text{left}(v)) \leq \alpha \cdot \lg \text{size}(\text{left}(v))$ . Combining these two inequalities with the recursive definition of height, we get

$$\alpha \cdot \lg \text{size}(v) < \text{height}(v) \leq \text{height}(\text{left}(v)) + 1 \leq \alpha \cdot \lg \text{size}(\text{left}(v)) + 1.$$

After some algebra, this simplifies to  $\text{size}(\text{left}(v)) > \text{size}(v)/2^{1/\alpha}$ . Combining this with the identity  $\text{size}(v) = \text{size}(\text{left}(v)) + \text{size}(\text{right}(v)) + 1$  and doing some more algebra gives us the inequality

$$\text{size}(\text{right}(v)) < (1 - 1/2^{1/\alpha})\text{size}(v) - 1.$$

Finally, we combine these two inequalities using the recursive definition of imbalance.

$$I(v) \geq \text{size}(\text{left}(v)) - \text{size}(\text{right}(v)) - 1 > (2/2^{1/\alpha} - 1)\text{size}(v)$$

Since  $\alpha$  is a constant bigger than 1, the factor in parentheses is a positive constant.  $\square$

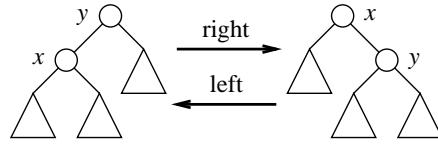
## 15.4 Scapegoat Trees

Finally, to handle both insertions and deletions efficiently, *scapegoat trees* use both of the previous techniques. We use partial rebuilding to re-balance the tree after insertions, and global rebuilding to re-balance the tree after deletions. Each search takes  $O(\log n)$  time in the worst case, and the amortized time for any insertion or deletion is also  $O(\log n)$ . There are a few small technical details left (which I won't describe), but no new ideas are required.

Once we've done the analysis, we can actually simplify the data structure. It's not hard to prove that at most one subtree (the *scapegoat*) is rebuilt during any insertion. Less obviously, we can even get the same amortized time bounds (except for a small constant factor) if we only maintain the three integers in addition to the actual tree: the size of the entire tree, the height of the entire tree, and the number of marked nodes. Whenever an insertion causes the tree to become unbalanced, we can compute the sizes of all the subtrees on the search path, starting at the new leaf and stopping at the scapegoat, in time proportional to the size of the scapegoat subtree. Since we need that much time to re-balance the scapegoat subtree, this computation increases the running time by only a small constant factor! Thus, unlike almost every other kind of balanced trees, scapegoat trees require only  $O(1)$  extra space.

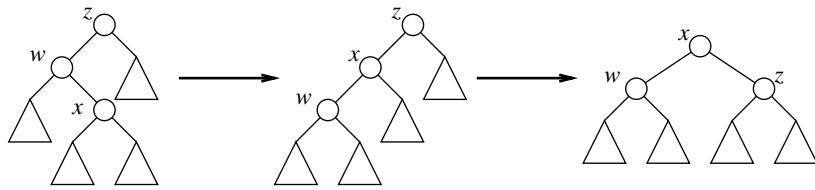
## 15.5 Rotations, Double Rotations, and Splaying

Another method for maintaining balance in binary search trees is by adjusting the shape of the tree locally, using an operation called a *rotation*. A rotation at a node  $x$  decreases its depth by one and increases its parent's depth by one. Rotations can be performed in constant time, since they only involve simple pointer manipulation.

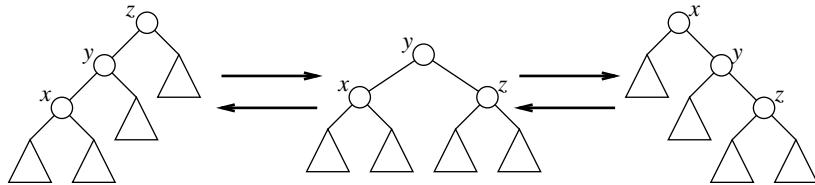


**Figure 1.** A right rotation at  $x$  and a left rotation at  $y$  are inverses.

For technical reasons, we will need to use rotations two at a time. There are two types of double rotations, which might be called *zig-zag* and *roller-coaster*. A zig-zag at  $x$  consists of two rotations at  $x$ , in opposite directions. A roller-coaster at a node  $x$  consists of a rotation at  $x$ 's parent followed by a rotation at  $x$ , both in the same direction. Each double rotation decreases the depth of  $x$  by two, leaves the depth of its parent unchanged, and increases the depth of its grandparent by either one or two, depending on the type of double rotation. Either type of double rotation can be performed in constant time.



**Figure 2.** A zig-zag at  $x$ . The symmetric case is not shown.



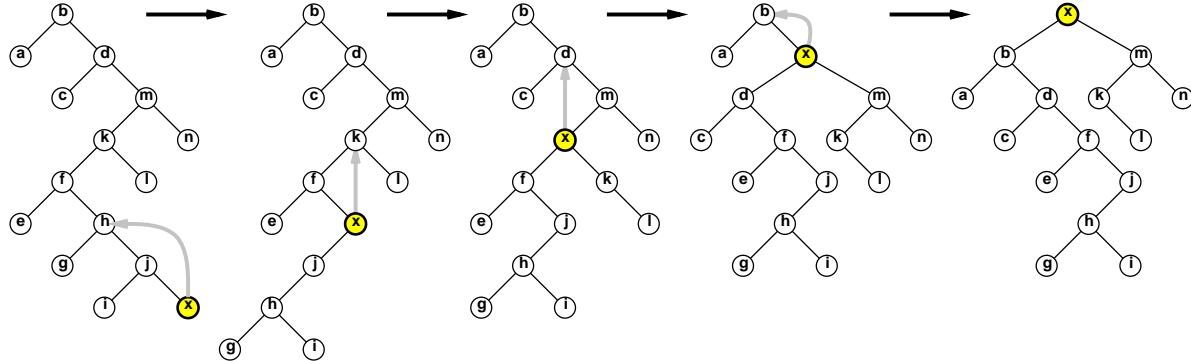
**Figure 3.** A right roller-coaster at  $x$  and a left roller-coaster at  $z$ .

Finally, a *splay* operation moves an arbitrary node in the tree up to the root through a series of double rotations, possibly with one single rotation at the end. Splaying a node  $v$  requires time proportional to  $\text{depth}(v)$ . (Obviously, this means the depth *before* splaying, since after splaying  $v$  is the root and thus has depth zero!)

## 15.6 Splay Trees

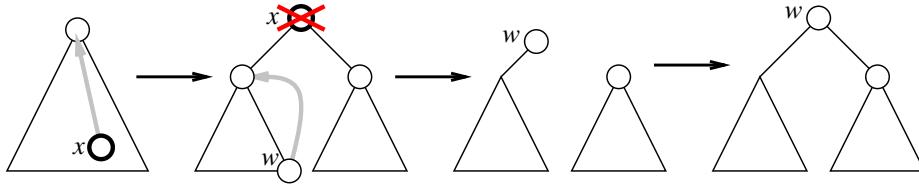
A *splay tree* is a binary search tree that is kept more or less balanced by splaying. Intuitively, after we access any node, we move it to the root with a splay operation. In more detail:

- **Search:** Find the node containing the key using the usual algorithm, or its predecessor or successor if the key is not present. Splay whichever node was found.



**Figure 4.** Splaying a node. Irrelevant subtrees are omitted for clarity.

- **Insert:** Insert a new node using the usual algorithm, then splay that node.
- **Delete:** Find the node  $x$  to be deleted, splay it, and then delete it. This splits the tree into two subtrees, one with keys less than  $x$ , the other with keys bigger than  $x$ . Find the node  $w$  in the left subtree with the largest key (the inorder predecessor of  $x$  in the original tree), splay it, and finally join it to the right subtree.



**Figure 5.** Deleting a node in a splay tree.

Each search, insertion, or deletion consists of a constant number of operations of the form *walk down to a node, and then splay it up to the root*. Since the walk down is clearly cheaper than the splay up, all we need to get good amortized bounds for splay trees is to derive good amortized bounds for a single splay.

Believe it or not, the easiest way to do this uses the potential method. We define the *rank* of a node  $v$  to be  $\lfloor \lg \text{size}(v) \rfloor$ , and the *potential* of a splay tree to be the sum of the ranks of its nodes:

$$\Phi := \sum_v \text{rank}(v) = \sum_v \lfloor \lg \text{size}(v) \rfloor$$

It's not hard to observe that a perfectly balanced binary tree has potential  $\Theta(n)$ , and a linear chain of nodes (a perfectly *unbalanced* tree) has potential  $\Theta(n \log n)$ .

The amortized analysis of splay trees boils down to the following lemma. Here,  $\text{rank}(v)$  denotes the rank of  $v$  before a (single or double) rotation, and  $\text{rank}'(v)$  denotes its rank afterwards. Recall that the amortized cost is defined to be the number of rotations plus the drop in potential.

**The Access Lemma.** *The amortized cost of a single rotation at any node  $v$  is at most  $1 + 3\text{rank}'(v) - 3\text{rank}(v)$ , and the amortized cost of a double rotation at any node  $v$  is at most  $3\text{rank}'(v) - 3\text{rank}(v)$ .*

Proving this lemma is a straightforward but tedious case analysis of the different types of rotations. For the sake of completeness, I'll give a proof (of a generalized version) in the next section.

By adding up the amortized costs of all the rotations, we find that the total amortized cost of splaying a node  $v$  is at most  $1 + 3\text{rank}'(v) - 3\text{rank}(v)$ , where  $\text{rank}'(v)$  is the rank of  $v$  after the entire splay. (The intermediate ranks cancel out in a nice telescoping sum.) But after the splay,  $v$  is the root! Thus,  $\text{rank}'(v) = \lfloor \lg n \rfloor$ , which implies that the amortized cost of a splay is at most  $3\lg n - 1 = O(\log n)$ .

We conclude that every insertion, deletion, or search in a splay tree takes  $O(\log n)$  amortized time.

### \*15.7 Other Optimality Properties

In fact, splay trees are optimal in several other senses. Some of these optimality properties follow easily from the following generalization of the Access Lemma.

Let's arbitrarily assign each node  $v$  a non-negative real *weight*  $w(v)$ . These weights are not actually stored in the splay tree, nor do they affect the splay algorithm in any way; they are only used to help with the analysis. We then redefine the *size*  $s(v)$  of a node  $v$  to be the sum of the weights of the descendants of  $v$ , including  $v$  itself:

$$s(v) := w(v) + s(\text{right}(v)) + s(\text{left}(v)).$$

If  $w(v) = 1$  for every node  $v$ , then the size of a node is just the number of nodes in its subtree, as in the previous section. As before, we define the *rank* of any node  $v$  to be  $r(v) = \lg s(v)$ , and the *potential* of a splay tree to be the sum of the ranks of all its nodes:

$$\Phi = \sum_v r(v) = \sum_v \lg s(v)$$

In the following lemma,  $r(v)$  denotes the rank of  $v$  before a (single or double) rotation, and  $r'(v)$  denotes its rank afterwards.

**The Generalized Access Lemma.** *For any assignment of non-negative weights to the nodes, the amortized cost of a single rotation at any node  $x$  is at most  $1 + 3r'(x) - 3r(x)$ , and the amortized cost of a double rotation at any node  $v$  is at most  $3r'(x) - 3r(x)$ .*

**Proof:** First consider a single rotation, as shown in Figure 1.

$$\begin{aligned} 1 + \Phi' - \Phi &= 1 + r'(x) + r'(y) - r(x) - r(y) && [\text{only } x \text{ and } y \text{ change rank}] \\ &\leq 1 + r'(x) - r(x) && [r'(y) \leq r(y)] \\ &\leq 1 + 3r'(x) - 3r(x) && [r'(x) \geq r(x)] \end{aligned}$$

Now consider a zig-zag, as shown in Figure 2. Only  $w$ ,  $x$ , and  $z$  change rank.

$$\begin{aligned} 2 + \Phi' - \Phi &= 2 + r'(w) + r'(x) + r'(z) - r(w) - r(x) - r(z) && [\text{only } w, x, z \text{ change rank}] \\ &\leq 2 + r'(w) + r'(x) + r'(z) - 2r(x) && [r(x) \leq r(w) \text{ and } r'(x) = r(z)] \\ &= 2 + (r'(w) - r'(x)) + (r'(z) - r'(x)) + 2(r'(x) - r(x)) \\ &= 2 + \lg \frac{s'(w)}{s'(x)} + \lg \frac{s'(z)}{s'(x)} + 2(r'(x) - r(x)) \\ &\leq 2 + 2 \lg \frac{s'(x)/2}{s'(x)} + 2(r'(x) - r(x)) && [s'(w) + s'(z) \leq s'(x), \lg \text{ is concave}] \\ &= 2(r'(x) - r(x)) \\ &\leq 3(r'(x) - r(x)) && [r'(x) \geq r(x)] \end{aligned}$$

Finally, consider a roller-coaster, as shown in Figure 3. Only  $x$ ,  $y$ , and  $z$  change rank.

$$\begin{aligned}
 2 + \Phi' - \Phi &= 2 + r'(x) + r'(y) + r'(z) - r(x) - r(y) - r(z) && [\text{only } x, y, z \text{ change rank}] \\
 &\leq 2 + r'(x) + r'(z) - 2r(x) && [r'(y) \leq r(z) \text{ and } r(x) \geq r(y)] \\
 &= 2 + (r(x) - r'(x)) + (r'(z) - r'(x)) + 3(r'(x) - r(x)) \\
 &= 2 + \lg \frac{s(x)}{s'(x)} + \lg \frac{s'(z)}{s'(x)} + 3(r'(x) - r(x)) \\
 &\leq 2 + 2\lg \frac{s'(x)/2}{s'(x)} + 3(r'(x) - r(x)) && [s(x) + s'(z) \leq s'(x), \lg \text{ is concave}] \\
 &= 3(r'(x) - r(x))
 \end{aligned}$$

This completes the proof.<sup>5</sup> □

Observe that this argument works for *arbitrary* non-negative vertex weights. By adding up the amortized costs of all the rotations, we find that the total amortized cost of splaying a node  $x$  is at most  $1 + 3r(\text{root}) - 3r(x)$ . (The intermediate ranks cancel out in a nice telescoping sum.)

This analysis has several immediate corollaries. The first corollary is that the amortized search time in a splay tree is within a constant factor of the search time in the best possible *static* binary search tree. Thus, if some nodes are accessed more often than others, the standard splay algorithm *automatically* keeps those more frequent nodes closer to the root, at least most of the time.

**Static Optimality Theorem.** *Suppose each node  $x$  is accessed at least  $t(x)$  times, and let  $T = \sum_x t(x)$ . The amortized cost of accessing  $x$  is  $O(\log T - \log t(x))$ .*

**Proof:** Set  $w(x) = t(x)$  for each node  $x$ . □

For any nodes  $x$  and  $z$ , let  $\text{dist}(x, z)$  denote the *rank distance* between  $x$  and  $z$ , that is, the number of nodes  $y$  such that  $\text{key}(x) \leq \text{key}(y) \leq \text{key}(z)$  or  $\text{key}(x) \geq \text{key}(y) \geq \text{key}(z)$ . In particular,  $\text{dist}(x, x) = 1$  for all  $x$ .

**Static Finger Theorem.** *For any fixed node  $f$  ('the finger'), the amortized cost of accessing  $x$  is  $O(\lg \text{dist}(f, x))$ .*

**Proof:** Set  $w(x) = 1/\text{dist}(x, f)^2$  for each node  $x$ . Then  $s(\text{root}) \leq \sum_{i=1}^{\infty} 2/i^2 = \pi^2/3 = O(1)$ , and  $r(x) \geq \lg w(x) = -2\lg \text{dist}(f, x)$ . □

Here are a few more interesting properties of splay trees, which I'll state without proof.<sup>6</sup> The proofs of these properties (especially the dynamic finger theorem) are considerably more complicated than the amortized analysis presented above.

**Working Set Theorem [14].** *The amortized cost of accessing node  $x$  is  $O(\log D)$ , where  $D$  is the number of distinct items accessed since the last time  $x$  was accessed. (For the first access to  $x$ , we set  $D = n$ .)*

---

<sup>5</sup>This proof is essentially taken verbatim from the original Sleator and Tarjan paper. Another proof technique, which may be more accessible, involves maintaining  $\lfloor \lg s(v) \rfloor$  tokens on each node  $v$  and arguing about the changes in token distribution caused by each single or double rotation. But I haven't yet internalized this approach enough to include it here.

<sup>6</sup>This list and the following section are taken almost directly from Erik Demaine's lecture notes [5].

**Scanning Theorem [17].** *Splaying all nodes in a splay tree in order, starting from any initial tree, requires  $O(n)$  total rotations.*

**Dynamic Finger Theorem [7, 6].** *Immediately after accessing node  $y$ , the amortized cost of accessing node  $x$  is  $O(\lg \text{dist}(x, y))$ .*

### \*15.8 Splay Tree Conjectures

Splay trees are conjectured to have many interesting properties in addition to the optimality properties that have been proved; I'll describe just a few of the more important ones.

The *Deque Conjecture* [17] considers the cost of dynamically maintaining two fingers  $l$  and  $r$ , starting on the left and right ends of the tree. Suppose at each step, we can move one of these two fingers either one step left or one step right; in other words, we are using the splay tree as a doubly-ended queue. Sundar\* proved that the total cost of  $m$  deque operations on an  $n$ -node splay tree is  $O((m+n)\alpha(m+n))$  [15]. More recently, Pettie later improved this bound to  $O(m\alpha^*(n))$  [13]. The Deque Conjecture states that the total cost is actually  $O(m+n)$ .

The *Traversal Conjecture* [14] states that accessing the nodes in a splay tree, in the order specified by a *preorder* traversal of any other binary tree with the same keys, takes  $O(n)$  time. This is generalization of the Scanning Theorem.

The *Unified Conjecture* [12] states that the time to access node  $x$  is  $O(\lg \min_y(D(y) + d(x, y)))$ , where  $D(y)$  is the number of *distinct* nodes accessed since the last time  $y$  was accessed. This would immediately imply both the Dynamic Finger Theorem, which is about spatial locality, and the Working Set Theorem, which is about temporal locality. Two other structures are known that satisfy the unified bound [4, 12].

Finally, the most important conjecture about splay trees, and one of the most important open problems about data structures, is that they are *dynamically optimal* [14]. Specifically, the cost of any sequence of accesses to a splay tree is conjectured to be at most a constant factor more than the cost of the best possible dynamic binary search tree *that knows the entire access sequence in advance*. To make the rules concrete, we consider binary search trees that can undergo *arbitrary* rotations after a search; the cost of a search is the number of key comparisons plus the number of rotations. We do not require that the rotations be on or even near the search path. This is an extremely strong conjecture!

No dynamically optimal binary search tree is known, even in the offline setting. However, three very similar  $O(\log \log n)$ -competitive binary search trees have been discovered in the last few years: *Tango trees* [8], *multisplay trees* [18], and *chain-splay trees* [11]. A recently-published geometric formulation of dynamic binary search trees [9] also offers significant hope for future progress.

## References

- [1] A. Andersson\*. Improving partial rebuilding by using simple balance criteria. *Proc. Workshop on Algorithms and Data Structures*, 393–402, 1989. Lecture Notes Comput. Sci. 382, Springer-Verlag.
- [2] A. Andersson. General balanced trees. *J. Algorithms* 30:1–28, 1999.
- [3] J. L. Bentley and J. B. Saxe\*. Decomposable searching problems I: Static-to-dynamic transformation. *J. Algorithms* 1(4):301–358, 1980.
- [4] M. Bădiu\* and E. D. Demaine. A simplified and dynamic unified structure. *Proc. 6th Latin American Symp. Theoretical Informatics*, 466–473, 2004. Lecture Notes Comput. Sci. 2976, Springer-Verlag.
- [5] J. Cohen\* and E. Demaine. 6.897: Advanced data structures (Spring 2005), Lecture 3, February 8 2005. <http://theory.csail.mit.edu/classes/6.897/spring05/lec.html>.
- [6] R. Cole. On the dynamic finger conjecture for splay trees. Part II: The proof. *SIAM J. Comput.* 30(1):44–85, 2000.

- [7] R. Cole, B. Mishra, J. Schmidt, and A. Siegel. On the dynamic finger conjecture for splay trees. Part I: Splay sorting  $\log n$ -block sequences. *SIAM J. Comput.* 30(1):1–43, 2000.
- [8] E. D. Demaine, D. Harmon\*, J. Iacono, and M. Pătrașcu\*. Dynamic optimality—almost. *Proc. 45th Ann. IEEE Symp. Foundations Comput. Sci.*, 484–490, 2004.
- [9] E. D. Demaine, D. Harmon\*, J. Iacono, D. Kane\*, and M. Pătrașcu\*. The geometry of binary search trees. *Proc. 20th ACM-SIAM Symp. Discrete Algorithms.*, 496–505, 2009.
- [10] I. Galperin\* and R. R. Rivest. Scapegoat trees. *Proc. 4th Ann. ACM-SIAM Symp. Discrete Algorithms*, 165–174, 1993.
- [11] G. Georgakopolous. Chain-splay trees, or, how to achieve and prove  $\log \log N$ -competitiveness by splaying. *Inform. Proc. Lett.* 106:37–34, 2008.
- [12] J. Iacono\*. Alternatives to splay trees with  $O(\log n)$  worst-case access times. *Proc. 12th Ann. ACM-SIAM Symp. Discrete Algorithms*, 516–522, 2001.
- [13] S. Pettie. Splay trees, Davenport-Schinzel sequences, and the deque conjecture. *Proc. 19th Ann. ACM-SIAM Symp. Discrete Algorithms*, 1115–1124, 2008.
- [14] D. D. Sleator and R. E. Tarjan. Self-adjusting binary search trees. *J. ACM* 32(3):652–686, 1985.
- [15] R. Sundar\*. On the Deque conjecture for the splay algorithm. *Combinatorica* 12(1):95–124, 1992.
- [16] R. E. Tarjan. *Data Structures and Network Algorithms*. CBMS-NSF Regional Conference Series in Applied Mathematics 44. SIAM, 1983.
- [17] R. E. Tarjan. Sequential access in splay trees takes linear time. *Combinatorica* 5(5):367–378, 1985.
- [18] C. C. Wang\*, J. Derryberry\*, and D. D. Sleator.  $O(\log \log n)$ -competitive dynamic binary search trees. *Proc. 17th Ann. ACM-SIAM Symp. Discrete Algorithms*, 374–383, 2006.

\*Starred authors were students at the time that the cited work was published.

## Exercises

1. (a) An  $n$ -node binary tree is *perfectly balanced* if either  $n \leq 1$ , or its two subtrees are perfectly balanced binary trees, each with at most  $\lfloor n/2 \rfloor$  nodes. Prove that  $I(v) = 0$  for every node  $v$  of any perfectly balanced tree.  
 (b) Prove that at most one subtree is rebalanced during a scapegoat tree insertion.
2. In a *dirty* binary search tree, each node is labeled either *clean* or *dirty*. The lazy deletion scheme used for scapegoat trees requires us to *purge* the search tree, keeping all the clean nodes and deleting all the dirty nodes, as soon as half the nodes become dirty. In addition, the purged tree should be perfectly balanced.
  - (a) Describe and analyze an algorithm to purge an arbitrary  $n$ -node dirty binary search tree in  $O(n)$  time. (Such an algorithm is necessary for scapegoat trees to achieve  $O(\log n)$  amortized insertion cost.)
  - \*(b) Modify your algorithm so that it uses only  $O(\log n)$  space, in addition to the tree itself. Don't forget to include the recursion stack in your space bound.
  - ★(c) Modify your algorithm so that it uses only  $O(1)$  additional space. In particular, your algorithm cannot call itself recursively at all.
3. Consider the following simpler alternative to splaying:

```
MOVEToRoot( $v$ ):
  while  $parent(v) \neq \text{NULL}$ 
    rotate at  $v$ 
```

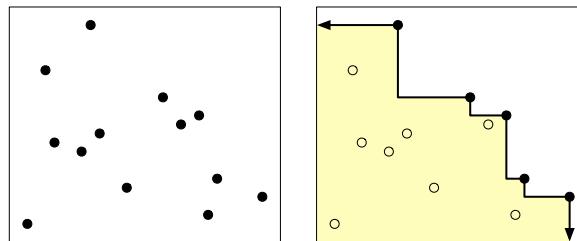
Prove that the amortized cost of `MoveToRoot` in an  $n$ -node binary tree can be  $\Omega(n)$ . That is, prove that for any integer  $k$ , there is a sequence of  $k$  `MoveToRoot` operations that require  $\Omega(kn)$  time to execute.

4. Suppose we want to maintain a dynamic set of values, subject to the following operations:

- `INSERT( $x$ )`: Add  $x$  to the set (if it isn't already there).
- `PRINT&DELETEBETWEEN( $a, b$ )`: Print every element  $x$  in the range  $a \leq x \leq b$ , in increasing order, and delete those elements from the set.

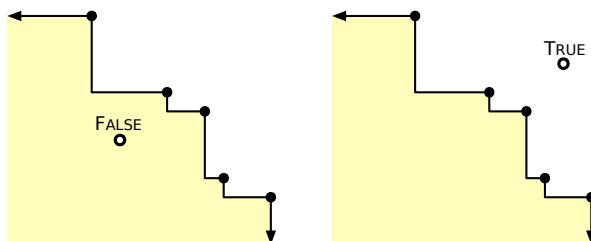
For example, if the current set is  $\{1, 5, 3, 4, 8\}$ , then `PRINT&DELETEBETWEEN(4, 6)` prints the numbers 4 and 5 and changes the set to  $\{1, 3, 8\}$ . For any underlying set, `PRINT&DELETEBETWEEN(-\infty, \infty)` prints its contents in increasing order and deletes everything.

- (a) Describe and analyze a data structure that supports these operations, each with amortized cost  $O(\log n)$ , where  $n$  is the maximum number of elements in the set.
- (b) What is the running time of your `INSERT` algorithm in the worst case?
- (c) What is the running time of your `PRINT&DELETEBETWEEN` algorithm in the worst case?
5. Let  $P$  be a set of  $n$  points in the plane. The *staircase* of  $P$  is the set of all points in the plane that have at least one point in  $P$  both above and to the right.



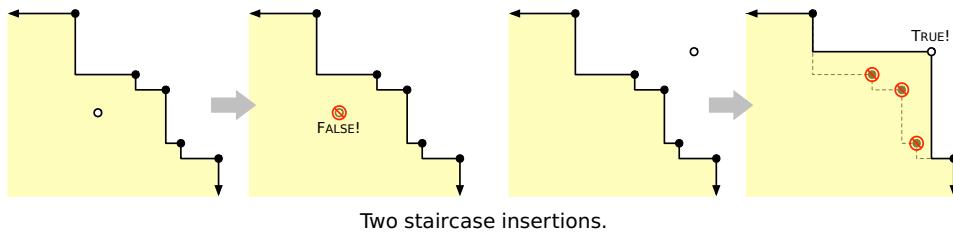
A set of points in the plane and its staircase (shaded).

- (a) Describe an algorithm to compute the staircase of a set of  $n$  points in  $O(n \log n)$  time.
- (b) Describe and analyze a data structure that stores the staircase of a set of points, and an algorithm `ABOVE?( $x, y$ )` that returns TRUE if the point  $(x, y)$  is above the staircase, or FALSE otherwise. Your data structure should use  $O(n)$  space, and your `ABOVE?` algorithm should run in  $O(\log n)$  time.



Two staircase queries.

- (c) Describe and analyze a data structure that maintains a staircase as new points are inserted. Specifically, your data structure should support a function  $\text{INSERT}(x, y)$  that adds the point  $(x, y)$  to the underlying point set and returns TRUE or FALSE to indicate whether the staircase of the set has changed. Your data structure should use  $O(n)$  space, and your  $\text{INSERT}$  algorithm should run in  $O(\log n)$  amortized time.



6. Say that a binary search tree is *augmented* if every node  $v$  also stores  $\text{size}(v)$ , the number of nodes in the subtree rooted at  $v$ .
- Show that a rotation in an augmented binary tree can be performed in constant time.
  - Describe an algorithm  $\text{SCAPEGOATSELECT}(k)$  that selects the  $k$ th smallest item in an augmented scapegoat tree in  $O(\log n)$  worst-case time. (The scapegoat trees presented in these notes are already augmented.)
  - Describe an algorithm  $\text{SPLAYSELECT}(k)$  that selects the  $k$ th smallest item in an augmented splay tree in  $O(\log n)$  amortized time.
  - Describe an algorithm  $\text{TREAPSELECT}(k)$  that selects the  $k$ th smallest item in an augmented treap in  $O(\log n)$  expected time.

7. Many applications of binary search trees attach a *secondary data structure* to each node in the tree, to allow for more complicated searches. Let  $T$  be an arbitrary binary tree. The secondary data structure at any node  $v$  stores exactly the same set of items as the subtree of  $T$  rooted at  $v$ . This secondary structure has size  $O(\text{size}(v))$  and can be built in  $O(\text{size}(v))$  time, where  $\text{size}(v)$  denotes the number of descendants of  $v$ .

The primary and secondary data structures are typically defined by different attributes of the data being stored. For example, to store a set of points in the plane, we could define the primary tree  $T$  in terms of the  $x$ -coordinates of the points, and define the secondary data structures in terms of their  $y$ -coordinate.

Maintaining these secondary structures complicates algorithms for keeping the top-level search tree balanced. Specifically, performing a rotation at any node  $v$  in the primary tree now requires  $O(\text{size}(v))$  time, because we have to rebuild one of the secondary structures (at the new child of  $v$ ). When we insert a new item into  $T$ , we must also insert into one or more secondary data structures.

- Overall, how much space does this data structure use *in the worst case*?
- How much space does this structure use if the primary search tree is perfectly balanced?
- Suppose the primary tree is a splay tree. Prove that the *amortized* cost of a splay (and therefore of a search, insertion, or deletion) is  $\Omega(n)$ . [Hint: This is easy!]

- (d) Now suppose the primary tree  $T$  is a scapegoat tree. How long does it take to rebuild the subtree of  $T$  rooted at some node  $v$ , as a function of  $\text{size}(v)$ ?
- (e) Suppose the primary tree and all the secondary trees are scapegoat trees. What is the amortized cost of a single insertion?
- \*(f) Finally, suppose the primary tree and every secondary tree is a treap. What is the worst-case expected time for a single insertion?
8. Let  $X = \langle x_1, x_2, \dots, x_m \rangle$  be a sequence of  $m$  integers, each from the set  $\{1, 2, \dots, n\}$ . We can visualize this sequence as a set of integer points in the plane, by interpreting each element  $x_i$  as the point  $(x_i, i)$ . The resulting point set, which we can also call  $X$ , has exactly one point on each row of the  $n \times m$  integer grid.
- (a) Let  $Y$  be an arbitrary set of integer points in the plane. Two points  $(x_1, y_1)$  and  $(x_2, y_2)$  in  $Y$  are *isolated* if (1)  $x_1 \neq x_2$  and  $y_1 \neq y_2$ , and (2) there is no other point  $(x, y) \in Y$  with  $x_1 \leq x \leq x_2$  and  $y_1 \leq y \leq y_2$ . If the set  $Y$  contains no isolated pairs of points, we call  $Y$  a *commune*.<sup>7</sup>
- Let  $X$  be an arbitrary set of points on the  $n \times n$  integer grid with exactly one point per row. Show that there is a commune  $Y$  that contains  $X$  and consists of  $O(n \log n)$  points.
- (b) Consider the following model of self-adjusting binary search trees. We interpret  $X$  as a sequence of accesses in a binary search tree. Let  $T_0$  denote the initial tree. In the  $i$ th round, we traverse the path from the root to node  $x_i$ , and then *arbitrarily reconfigure* some subtree  $S_i$  of the current search tree  $T_{i-1}$  to obtain the next search tree  $T_i$ . The only restriction is that the subtree  $S_i$  must contain both  $x_i$  and the root of  $T_{i-1}$ . (For example, in a splay tree,  $S_i$  is the search path to  $x_i$ .) The *cost* of the  $i$ th access is the number of nodes in the subtree  $S_i$ .
- Prove that the minimum cost of executing an access sequence  $X$  in this model is at least the size of the smallest commune containing the corresponding point set  $X$ . [Hint: Lowest common ancestor.]
- \*(c) Suppose  $X$  is a *random* permutation of the integers  $1, 2, \dots, n$ . Use the lower bound in part (b) to prove that the expected minimum cost of executing  $X$  is  $\Omega(n \log n)$ .
- ★(d) Describe a polynomial-time algorithm to compute (or even approximate up to constant factors) the smallest commune containing a given set  $X$  of integer points, with at most one point per row. Alternately, prove that the problem is NP-hard.

---

<sup>7</sup>Demaine et al. [9] refer to communes as *arborally satisfied sets*.

*E pluribus unum (Out of many, one)*

— Official motto of the United States of America

**John:** Who's your daddy? C'mon, you know who your daddy is! Who's your daddy?  
D'Argo, tell him who his daddy is!"

**D'Argo:** I'm your daddy.

— *Farscape*, "Thanks for Sharing" (June 15, 2001)

What rolls down stairs, alone or in pairs, rolls over your neighbor's dog?

What's great for a snack, and fits on your back? It's Log, Log, Log!

It's Log! It's Log! It's big, it's heavy, it's wood!  
It's Log! It's Log! It's better than bad, it's good!

— Ren & Stimpy, "Stimpy's Big Day/The Big Shot" (August 11, 1991)  
lyrics by John Kricfalusi

The thing's hollow - it goes on forever - and - oh my God! - it's full of stars!

— Capt. David Bowman's last words(?)  
*2001: A Space Odyssey* by Arthur C. Clarke (1968)

## 16 Data Structures for Disjoint Sets

In this lecture, we describe some methods for maintaining a collection of disjoint sets. Each set is represented as a pointer-based data structure, with one node per element. We will refer to the elements as either ‘objects’ or ‘nodes’, depending on whether we want to emphasize the set abstraction or the actual data structure. Each set has a unique ‘leader’ element, which identifies the set. (Since the sets are always disjoint, the same object cannot be the leader of more than one set.) We want to support the following operations.

- **MAKESET( $x$ ):** Create a new set  $\{x\}$  containing the single element  $x$ . The object  $x$  must not appear in any other set in our collection. The leader of the new set is obviously  $x$ .
- **FIND( $x$ ):** Find (the leader of) the set containing  $x$ .
- **UNION( $A, B$ ):** Replace two sets  $A$  and  $B$  in our collection with their union  $A \cup B$ . For example,  $\text{UNION}(A, \text{MAKESET}(x))$  adds a new element  $x$  to an existing set  $A$ . The sets  $A$  and  $B$  are specified by arbitrary elements, so  $\text{UNION}(x, y)$  has exactly the same behavior as  $\text{UNION}(\text{FIND}(x), \text{FIND}(y))$ .

Disjoint set data structures have lots of applications. For instance, Kruskal’s minimum spanning tree algorithm relies on such a data structure to maintain the components of the intermediate spanning forest. Another application is maintaining the connected components of a graph as new vertices and edges are added. In both these applications, we can use a disjoint-set data structure, where we maintain a set for each connected component, containing that component’s vertices.

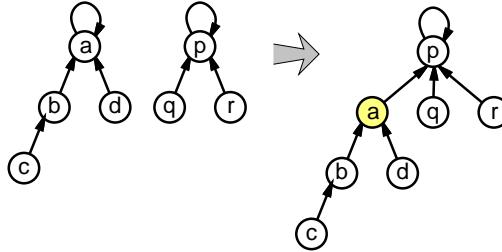
### 16.1 Reversed Trees

One of the easiest ways to store sets is using trees, in which each node represents a single element of the set. Each node points to another node, called its *parent*, except for the leader of each set, which points to itself and thus is the root of the tree. **MAKESET** is trivial. **FIND** traverses parent pointers up to the leader. **UNION** just redirects the parent pointer of one leader to the other. Unlike most tree data structures, nodes do *not* have pointers down to their children.

```
MAKESET( $x$ ):
  parent( $x$ )  $\leftarrow x$ 
```

```
FIND( $x$ ):
  while  $x \neq \text{parent}(x)$ 
     $x \leftarrow \text{parent}(x)$ 
  return  $x$ 
```

```
UNION( $x, y$ ):
   $\bar{x} \leftarrow \text{FIND}(x)$ 
   $\bar{y} \leftarrow \text{FIND}(y)$ 
  parent( $\bar{y}$ )  $\leftarrow \bar{x}$ 
```



Merging two sets stored as trees. Arrows point to parents. The shaded node has a new parent.

MAKE-SET clearly takes  $\Theta(1)$  time, and UNION requires only  $O(1)$  time in addition to the two FINDs. The running time of  $\text{FIND}(x)$  is proportional to the depth of  $x$  in the tree. It is not hard to come up with a sequence of operations that results in a tree that is a long chain of nodes, so that FIND takes  $\Theta(n)$  time in the worst case.

However, there is an easy change we can make to our UNION algorithm, called *union by depth*, so that the trees always have logarithmic depth. Whenever we need to merge two trees, we always make the root of the *shallow*er tree a child of the *deep*er one. This requires us to also maintain the depth of each tree, but this is quite easy.

```
MAKESET( $x$ ):
  parent( $x$ )  $\leftarrow x$ 
  depth( $x$ )  $\leftarrow 0$ 
```

```
FIND( $x$ ):
  while  $x \neq \text{parent}(x)$ 
     $x \leftarrow \text{parent}(x)$ 
  return  $x$ 
```

```
UNION( $x, y$ )
   $\bar{x} \leftarrow \text{FIND}(x)$ 
   $\bar{y} \leftarrow \text{FIND}(y)$ 
  if  $\text{depth}(\bar{x}) > \text{depth}(\bar{y})$ 
    parent( $\bar{y}$ )  $\leftarrow \bar{x}$ 
  else
    parent( $\bar{x}$ )  $\leftarrow \bar{y}$ 
    if  $\text{depth}(\bar{x}) = \text{depth}(\bar{y})$ 
      depth( $\bar{y}$ )  $\leftarrow \text{depth}(\bar{y}) + 1$ 
```

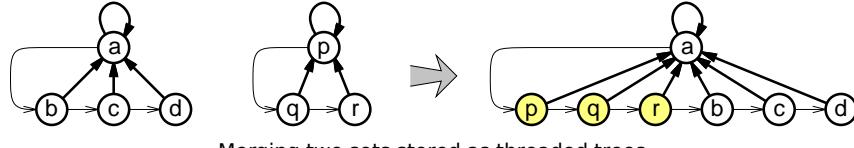
With this new rule in place, it's not hard to prove by induction that for any set leader  $\bar{x}$ , the size of  $\bar{x}$ 's set is at least  $2^{\text{depth}(\bar{x})}$ , as follows. If  $\text{depth}(\bar{x}) = 0$ , then  $\bar{x}$  is the leader of a singleton set. For any  $d > 0$ , when  $\text{depth}(\bar{x})$  becomes  $d$  for the first time,  $\bar{x}$  is becoming the leader of the union of two sets, both of whose leaders had depth  $d - 1$ . By the inductive hypothesis, both component sets had at least  $2^{d-1}$  elements, so the new set has at least  $2^d$  elements. Later UNION operations might add elements to  $\bar{x}$ 's set without changing its depth, but that only helps us.

Since there are only  $n$  elements altogether, the maximum depth of any set is  $\lg n$ . We conclude that if we use union by depth, both FIND and UNION run in  $\Theta(\log n)$  time in the worst case.

## 16.2 Shallow Threaded Trees

Alternately, we could just have every object keep a pointer to the leader of its set. Thus, each set is represented by a shallow tree, where the leader is the root and all the other elements are its children. With this representation, MAKESET and FIND are completely trivial. Both operations clearly run in constant time. UNION is a little more difficult, but not much. Our algorithm sets all the leader pointers in one set to point to the leader of the other set. To do this, we need a method to visit every element

in a set; we will ‘thread’ a linked list through each set, starting at the set’s leader. The two threads are merged in the UNION algorithm in constant time.



Merging two sets stored as threaded trees.

Bold arrows point to leaders; lighter arrows form the threads. Shaded nodes have a new leader.

MAKESET( $x$ ):  
leader( $x$ )  $\leftarrow x$   
next( $x$ )  $\leftarrow x$

FIND( $x$ ):  
return leader( $x$ )

UNION( $x, y$ ):  
 $\bar{x} \leftarrow \text{FIND}(x)$   
 $\bar{y} \leftarrow \text{FIND}(y)$   
 $y \leftarrow \bar{y}$   
leader( $y$ )  $\leftarrow \bar{x}$   
while (next( $y$ )  $\neq \text{NULL}$ )  
     $y \leftarrow \text{next}(y)$   
    leader( $y$ )  $\leftarrow \bar{x}$   
next( $y$ )  $\leftarrow \text{next}(\bar{x})$   
next( $\bar{x}$ )  $\leftarrow \bar{y}$

The worst-case running time of UNION is a constant times the size of the *larger* set. Thus, if we merge a one-element set with another  $n$ -element set, the running time can be  $\Theta(n)$ . Generalizing this idea, it is quite easy to come up with a sequence of  $n$  MAKESET and  $n - 1$  UNION operations that requires  $\Theta(n^2)$  time to create the set  $\{1, 2, \dots, n\}$  from scratch.

WORSTCASESEQUENCE( $n$ ):  
MAKESET(1)  
for  $i \leftarrow 2$  to  $n$   
    MAKESET( $i$ )  
    UNION(1,  $i$ )

We are being stupid in two different ways here. One is the order of operations in WORSTCASESEQUENCE. Obviously, it would be more efficient to merge the sets in the other order, or to use some sort of divide and conquer approach. Unfortunately, we can’t fix this; we don’t get to decide how our data structures are used! The other is that we always update the leader pointers in the larger set. To fix this, we add a comparison inside the UNION algorithm to determine which set is smaller. This requires us to maintain the size of each set, but that’s easy.

MAKEWEIGHTEDSET( $x$ ):  
leader( $x$ )  $\leftarrow x$   
next( $x$ )  $\leftarrow x$   
size( $x$ )  $\leftarrow 1$

WEIGHTEDUNION( $x, y$ ):  
 $\bar{x} \leftarrow \text{FIND}(x)$   
 $\bar{y} \leftarrow \text{FIND}(y)$   
if size( $\bar{x}$ )  $>$  size( $\bar{y}$ )  
    UNION( $\bar{x}, \bar{y}$ )  
    size( $\bar{x}$ )  $\leftarrow \text{size}(\bar{x}) + \text{size}(\bar{y})$   
else  
    UNION( $\bar{y}, \bar{x}$ )  
    size( $\bar{x}$ )  $\leftarrow \text{size}(\bar{x}) + \text{size}(\bar{y})$

The new WEIGHTEDUNION algorithm still takes  $\Theta(n)$  time to merge two  $n$ -element sets. However, in an amortized sense, this algorithm is much more efficient. Intuitively, before we can merge two large sets, we have to perform a large number of MAKEWEIGHTEDSET operations.

**Theorem 1.** A sequence of  $m$  MAKEWEIGHTEDSET operations and  $n$  WEIGHTEDUNION operations takes  $O(m + n \log n)$  time in the worst case.

**Proof:** Whenever the leader of an object  $x$  is changed by a WEIGHTEDUNION, the size of the set containing  $x$  increases by at least a factor of two. By induction, if the leader of  $x$  has changed  $k$  times, the set containing  $x$  has at least  $2^k$  members. After the sequence ends, the largest set contains at most  $n$  members. (Why?) Thus, the leader of any object  $x$  has changed at most  $\lfloor \lg n \rfloor$  times.

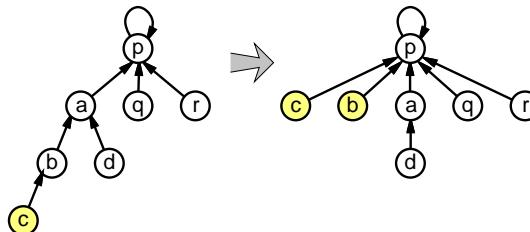
Since each WEIGHTEDUNION reduces the number of sets by one, there are  $m - n$  sets at the end of the sequence, and at most  $n$  objects are *not* in singleton sets. Since each of the non-singleton objects had  $O(\log n)$  leader changes, the total amount of work done in updating the leader pointers is  $O(n \log n)$ .  $\square$

The aggregate method now implies that each WEIGHTEDUNION has *amortized cost*  $O(\log n)$ .

### 16.3 Path Compression

Using unthreaded trees, FIND takes logarithmic time and everything else is constant; using threaded trees, UNION takes logarithmic amortized time and everything else is constant. A third method allows us to get both of these operations to have *almost* constant running time.

We start with the original unthreaded tree representation, where every object points to a parent. The key observation is that in any FIND operation, once we determine the leader of an object  $x$ , we can speed up future FINDs by redirecting  $x$ 's parent pointer directly to that leader. In fact, we can change the parent pointers of all the ancestors of  $x$  all the way up to the root; this is easiest if we use recursion for the initial traversal up the tree. This modification to FIND is called *path compression*.



Path compression during  $\text{Find}(c)$ . Shaded nodes have a new parent.

```

FIND( $x$ )
  if  $x \neq \text{parent}(x)$ 
     $\text{parent}(x) \leftarrow \text{FIND}(\text{parent}(x))$ 
  return  $\text{parent}(x)$ 

```

If we use path compression, the ‘depth’ field we used earlier to keep the trees shallow is no longer correct, and correcting it would take way too long. But this information still ensures that FIND runs in  $\Theta(\log n)$  time in the worst case, so we’ll just give it another name: *rank*. The following algorithm is usually called *union by rank*:

```

MAKESET( $x$ ):
   $\text{parent}(x) \leftarrow x$ 
   $\text{rank}(x) \leftarrow 0$ 

```

```

UNION( $x, y$ )
   $\bar{x} \leftarrow \text{FIND}(x)$ 
   $\bar{y} \leftarrow \text{FIND}(y)$ 
  if  $\text{rank}(\bar{x}) > \text{rank}(\bar{y})$ 
     $\text{parent}(\bar{y}) \leftarrow \bar{x}$ 
  else
     $\text{parent}(\bar{x}) \leftarrow \bar{y}$ 
    if  $\text{rank}(\bar{x}) = \text{rank}(\bar{y})$ 
       $\text{rank}(\bar{y}) \leftarrow \text{rank}(\bar{y}) + 1$ 

```

FIND still runs in  $O(\log n)$  time in the worst case; path compression increases the cost by only most a constant factor. But we have good reason to suspect that this upper bound is no longer tight. Our new algorithm memoizes the results of each FIND, so if we are asked to FIND the same item twice in a row, the second call returns in constant time. Splay trees used a similar strategy to achieve their optimal amortized cost, but our up-trees have fewer constraints on their structure than binary search trees, so we should get even better performance.

This intuition is exactly correct, but it takes a bit of work to define precisely *how* much better the performance is. As a first approximation, we will prove below that the amortized cost of a FIND operation is bounded by the *iterated logarithm* of  $n$ , denoted  $\log^* n$ , which is the number of times one must take the logarithm of  $n$  before the value is less than 1:

$$\lg^* n = \begin{cases} 1 & \text{if } n \leq 2, \\ 1 + \lg^*(\lg n) & \text{otherwise.} \end{cases}$$

Our proof relies on several useful properties of ranks, which follow directly from the UNION and FIND algorithms.

- If a node  $x$  is not a set leader, then the rank of  $x$  is smaller than the rank of its parent.
- Whenever  $\text{parent}(x)$  changes, the new parent has larger rank than the old parent.
- Whenever the leader of  $x$ 's set changes, the new leader has larger rank than the old leader.
- The size of any set is exponential in the rank of its leader:  $\text{size}(\bar{x}) \geq 2^{\text{rank}(\bar{x})}$ . (This is easy to prove by induction, hint, hint.)
- In particular, since there are only  $n$  objects, the highest possible rank is  $\lfloor \lg n \rfloor$ .
- For any integer  $r$ , there are at most  $n/2^r$  objects of rank  $r$ .

Only the last property requires a clever argument to prove. Fix your favorite integer  $r$ . Observe that only set leaders can change their rank. Whenever the rank of any set leader  $\bar{x}$  changes from  $r - 1$  to  $r$ , mark all the objects in  $\bar{x}$ 's set. Since leader ranks can only increase over time, each object is marked at most once. There are  $n$  objects altogether, and any object with rank  $r$  marks at least  $2^r$  objects. It follows that there are at most  $n/2^r$  objects with rank  $r$ , as claimed.

#### \*16.4 $O(\log^* n)$ Amortized Time

The following analysis of path compression was discovered just a few years ago by Raimund Seidel and Micha Sharir.<sup>1</sup> Previous proofs<sup>2</sup> relied on complicated charging schemes or potential-function arguments; Seidel and Sharir's analysis relies on a comparatively simple recursive decomposition.<sup>3</sup>

Seidel and Sharir phrase their analysis in terms of two more general operations on set forests. Their more general COMPRESS operation compresses *any* directed path, not just paths that lead to the root. The new SHATTER operation makes every node on a root-to-leaf path into its own parent.

---

<sup>1</sup>Raimund Seidel and Micha Sharir. Top-down analysis of path compression. *SIAM J. Computing* 34(3):515–525, 2005.

<sup>2</sup>Robert E. Tarjan. Efficiency of a good but not linear set union algorithm. *J. Assoc. Comput. Mach.* 22:215–225, 1975.

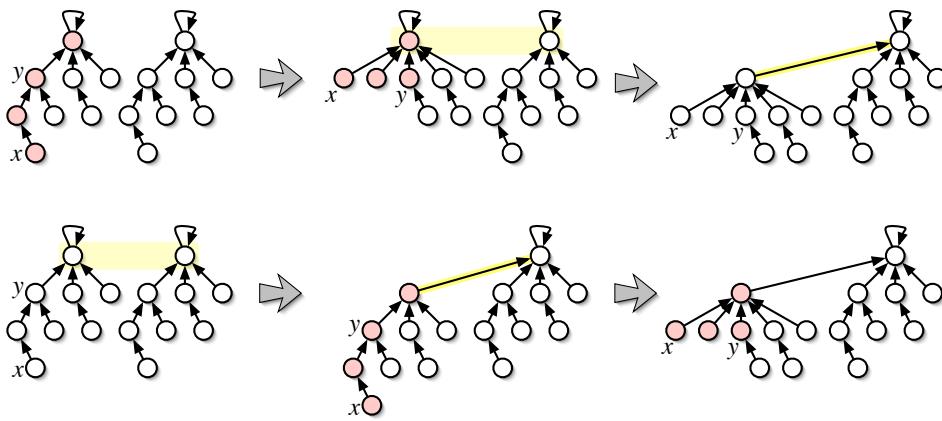
<sup>3</sup>For an even simpler analysis, see: Seth Pettie. Applications of forbidden 0-1 matrices to search tree and path compression-based data structures. *Proc. 21st Ann. ACM-SIAM Symp. Discrete Algorithms*, 1457–1467, 2010.

```
COMPRESS( $x, y$ ):
  « $y$  must be an ancestor of  $x$ »
  if  $x \neq y$ 
    COMPRESS( $\text{parent}(x), y$ )
     $\text{parent}(x) \leftarrow \text{parent}(y)$ 
```

```
SHATTER( $x$ ):
  if  $\text{parent}(x) \neq x$ 
    SHATTER( $\text{parent}(x)$ )
     $\text{parent}(x) \leftarrow x$ 
```

Clearly, the running time of FIND( $x$ ) operation is dominated by the running time of COMPRESS( $x, y$ ), where  $y$  is the leader of the set containing  $x$ . This implies that we can prove the upper bound by analyzing an arbitrary sequence of UNION and COMPRESS operations. Moreover, we can assume that the arguments to each UNION operation are set leaders, so that each UNION takes only constant worst-case time.

Finally, since each call to COMPRESS specifies the top node in the path to be compressed, we can reorder the sequence of operations, so that every UNION occurs before any COMPRESS, without changing the number of pointer assignments.



Top row: A COMPRESS followed by a UNION. Bottom row: The same operations in the opposite order.

Each UNION requires only constant time in the worst case, so we only need to analyze the amortized cost of COMPRESS. The running time of COMPRESS is proportional to the number of parent pointer assignments, plus  $O(1)$  overhead, so we will phrase our analysis in terms of pointer assignments. Let  $T(m, n, r)$  denote the worst case number of pointer assignments in any sequence of at most  $m$  COMPRESS operations, executed on a forest of at most  $n$  nodes, with maximum rank at most  $r$ .

The following trivial upper bound will be the base case for our recursive argument.

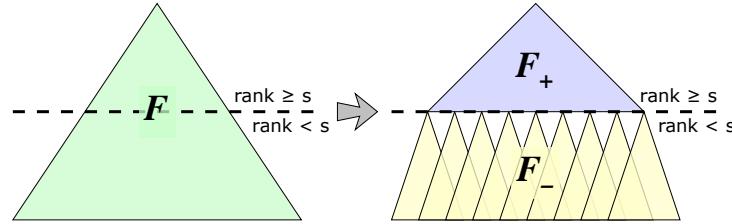
**Theorem 2.**  $T(m, n, r) \leq nr$

**Proof:** Each node can change parents at most  $r$  times, because each new parent has higher rank than the previous parent.  $\square$

Fix a forest  $F$  of  $n$  nodes with maximum rank  $r$ , and a sequence  $C$  of  $m$  COMPRESS operations on  $F$ , and let  $T(F, C)$  denote the total number of pointer assignments executed by this sequence.

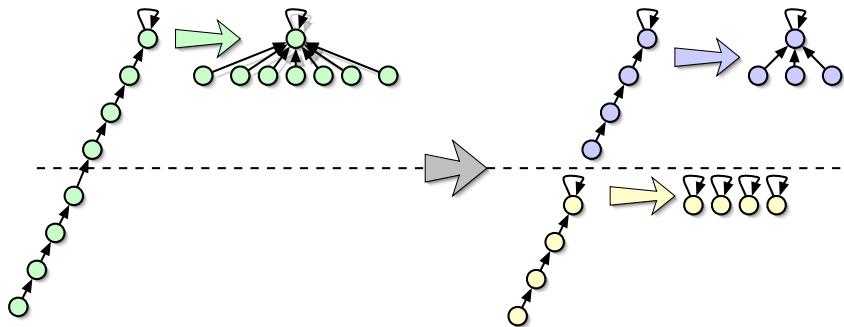
Let  $s$  be an arbitrary positive rank. Partition  $F$  into two sub-forests: a ‘low’ forest  $F_-$  containing all nodes with rank at most  $s$ , and a ‘high’ forest  $F_+$  containing all nodes with rank greater than  $s$ . Since ranks increase as we follow parent pointers, every ancestor of a high node is another high node. Let  $n_-$  and  $n_+$  denote the number of nodes in  $F_-$  and  $F_+$ , respectively. Finally, let  $m_+$  denote the number of COMPRESS operations that involve any node in  $F_+$ , and let  $m_- = m - m_+$ .

Any sequence of COMPRESS operations on  $F$  can be decomposed into a sequence of COMPRESS operations on  $F_+$ , plus a sequence of COMPRESS and SHATTER operations on  $F_-$ , with the same total



Splitting the forest  $F$  (in this case, a single tree) into sub-forests  $F_+$  and  $F_-$  at rank  $s$ .

cost. This requires only one small modification to the code: We forbid any low node from having a high parent. Specifically, if  $x$  is a low node and  $y$  is a high node, we replace any assignment  $\text{parent}(x) \leftarrow y$  with  $\text{parent}(x) \leftarrow x$ .



A Compress operation in  $F$  splits into a Compress operation in  $F_+$  and a Shatter operation in  $F_-$

This modification is equivalent to the following reduction:

```

COMPRESS( $x, y, F$ ):    $\langle\langle y \text{ is an ancestor of } x \rangle\rangle$ 
    if  $\text{rank}(x) > r$ 
        COMPRESS( $x, y, F_+$ )    $\langle\langle \text{in } C_+ \rangle\rangle$ 
    else if  $\text{rank}(y) \leq r$ 
        COMPRESS( $x, y, F_-$ )    $\langle\langle \text{in } C_- \rangle\rangle$ 
    else
         $z \leftarrow \text{highest ancestor of } x \text{ in } F \text{ with rank at most } r$ 
        COMPRESS( $\text{parent}_F(z), y, F_+$ )    $\langle\langle \text{in } C_+ \rangle\rangle$ 
        SHATTER( $x, z, F_-$ )
         $\text{parent}(z) \leftarrow z$    (*)
    
```

The pointer assignment in the last line looks redundant, but it is actually necessary for the analysis. Each execution of line (\*) mirrors an assignment of the form  $\text{parent}(x) \leftarrow y$ , where  $x$  is a low node,  $y$  is a high node, and the previous parent of  $x$  was a high node. Each of these ‘redundant’ assignments happens immediately after a COMPRESS in the top forest, so we perform at most  $m_+$  redundant assignments.

Each node  $x$  is touched by at most one SHATTER operation, so the total number of pointer reassignments in all the SHATTER operations is at most  $n$ .

Thus, by partitioning the forest  $F$  into  $F_+$  and  $F_-$ , we have also partitioned the sequence  $C$  of COMPRESS operations into subsequences  $C_+$  and  $C_-$ , with respective lengths  $m_+$  and  $m_-$ , such that the following inequality holds:

$$T(F, C) \leq T(F_+, C_+) + T(F_-, C_-) + m_+ + n$$

Since there are only  $n/2^i$  nodes of any rank  $i$ , we have  $n_+ \leq \sum_{i>s} n/2^i = n/2^s$ . The number of different ranks in  $F_+$  is  $r - s < r$ . Thus, Theorem 2 implies the upper bound

$$T(F_+, C_+) < rn/2^s.$$

Let us fix  $s = \lg r$ , so that  $T(F_+, C_+) \leq n$ . We can now simplify our earlier recurrence to

$$T(F, C) \leq T(F_-, C_-) + m_+ + 2n,$$

or equivalently,

$$T(F, C) - m \leq T(F_-, C_-) - m_- + 2n.$$

Since this argument applies to *any* forest  $F$  and *any* sequence  $C$ , we have just proved that

$$T'(m, n, r) \leq T'(m, n, \lfloor \lg r \rfloor) + 2n,$$

where  $T'(m, n, r) = T(m, n, r) - m$ . The solution to this recurrence is  $T'(n, m, r) \leq 2n \lg^* r$ . Voilá!

**Theorem 3.**  $T(m, n, r) \leq m + 2n \lg^* r$

### \*16.5 Turning the Crank

There is one place in the preceding analysis where we have significant room for improvement. Recall that we bounded the total cost of the operations on  $F_+$  using the trivial upper bound from Theorem 2. But we just proved a better upper bound in Theorem 3! We can apply precisely the same strategy, using Theorem 3 instead of Theorem 2, to improve the bound even more.

Suppose we fix  $s = \lg^* r$ , so that  $n_+ = n/2^{\lg^* r}$ . Theorem 3 implies that

$$T(F_+, C_+) \leq m_+ + 2n \frac{\lg^* r}{2^{\lg^* r}} \leq m_+ + 2n.$$

This implies the recurrence

$$T(F, C) \leq T(F_-, C_-) + 2m_+ + 3n,$$

which in turn implies that

$$T''(m, n, r) \leq T''(m, n, \lg^* r) + 3n,$$

where  $T''(m, n, r) = T(m, n, r) - 2m$ . The solution to this equation is  $T(m, n, r) \leq 2m + 3n \lg^{**} r$ , where  $\lg^{**} r$  is the *iterated* iterated logarithm of  $r$ :

$$\lg^{**} r = \begin{cases} 1 & \text{if } r \leq 2, \\ 1 + \lg^{**}(\lg^* r) & \text{otherwise.} \end{cases}$$

Naturally we can apply the same improvement strategy again, and again, as many times as we like, each time producing a tighter upper bound. Applying the reduction  $c$  times, for any positive integer  $c$ , gives us

$$T(m, n, r) \leq cm + (c+1)n \lg^{*c} r$$

where

$$\lg^{*c} r = \begin{cases} \lg r & \text{if } c = 0, \\ 1 & \text{if } r \leq 2, \\ 1 + \lg^{*c}(\lg^{*c-1} r) & \text{otherwise.} \end{cases}$$

Each time we ‘turn the crank’, the dependence on  $m$  increases, while the dependence on  $n$  and  $r$  decreases. For sufficiently large values of  $c$ , the  $cm$  term dominates the time bound, and further iterations only make things worse. The point of diminishing returns can be estimated by *the minimum number of stars* such that  $\lg^{*...*} r$  is smaller than a constant:

$$\alpha(r) = \min \left\{ c \geq 1 \mid \lg^{*^c} n \leq 3 \right\}.$$

(The threshold value 3 is used here because  $\lg^{*^c} 5 \geq 2$  for all  $c$ .) By setting  $c = \alpha(r)$ , we obtain our final upper bound.

**Theorem 4.**  $T(m, n, r) \leq m\alpha(r) + 3n(\alpha(r) + 1)$

We can assume without loss of generality that  $m \geq n$  by ignoring any singleton sets, so this upper bound can be further simplified to  $T(m, n, r) = O(m\alpha(r)) = O(m\alpha(n))$ . It follows that if we use union by rank, FIND with path compression runs in  **$O(\alpha(n))$  amortized time**.

Even this upper bound is somewhat conservative if  $m$  is larger than  $n$ . A closer estimate is given by the function

$$\alpha(m, n) = \min \left\{ c \geq 1 \mid \log^{*^c} (\lg n) \leq m/n \right\}.$$

It’s not hard to prove that if  $m = \Theta(n)$ , then  $\alpha(m, n) = \Theta(\alpha(n))$ . On the other hand, if  $m \geq n \lg^{****} n$ , for any constant number of stars, then  $\alpha(m, n) = O(1)$ . So even if the number of FIND operations is only *slightly* larger than the number of nodes, the amortized cost of each FIND is *constant*.

$O(\alpha(m, n))$  is actually a *tight* upper bound for the amortized cost of path compression; there are no more tricks that will improve the analysis further. More surprisingly, this is the best amortized bound we obtain for *any* pointer-based data structure for maintaining disjoint sets; the amortized cost of *every* FIND algorithm is at least  $\Omega(\alpha(m, n))$ . The proof of the matching lower bound is, unfortunately, far beyond the scope of this class.<sup>4</sup>

## 16.6 The Ackermann Function and its Inverse

The iterated logarithms that fell out of our analysis of path compression are the inverses of a hierarchy of recursive functions defined by Wilhelm Ackermann in 1928.<sup>5</sup>

$$2 \uparrow^c n = \begin{cases} 2 & \text{if } n = 1 \\ 2n & \text{if } c = 0 \\ 2 \uparrow^{c-1} (2 \uparrow^c (n-1)) & \text{otherwise} \end{cases}$$

For each fixed  $c$ , the function  $2 \uparrow^c n$  is monotonically increasing in  $n$ , and these functions grow *incredibly* faster as the index  $c$  increases.  $2 \uparrow n$  is the familiar power function  $2^n$ .  $2 \uparrow\uparrow n$  is the *tower* function  $2^{2^{\cdot^{\cdot^2}}}_n$ ; this function is also sometimes called *tetration*. John Conway named  $2 \uparrow\uparrow\uparrow n$  the *wower* function:  $2 \uparrow\uparrow\uparrow n = \underbrace{2 \uparrow\uparrow 2 \uparrow\uparrow \cdots \uparrow\uparrow 2}_n$ . And so on, *et cetera, ad infinitum*.

<sup>4</sup>Robert E. Tarjan. A class of algorithms which require non-linear time to maintain disjoint sets. *J. Comput. Syst. Sci.* 19:110–127, 1979.

<sup>5</sup>Ackermann didn’t define his functions this way—I’m actually describing a slightly cleaner hierarchy defined 35 years later by R. Creighton Buck—but the exact details of the definition are surprisingly irrelevant! The mnemonic up-arrow notation for these functions was introduced by Don Knuth in the 1970s.

For any fixed  $c$ , the function  $\log^{*c} n$  is the inverse of the function  $2 \uparrow^{c+1} n$ , the  $(c+1)$ th row in the Ackerman hierarchy. Thus, for any remotely reasonable values of  $n$ , say  $n \leq 2^{256}$ , we have  $\log^* n \leq 5$ ,  $\log^{**} n \leq 4$ , and  $\log^{*c} n \leq 3$  for any  $c \geq 3$ .

The function  $\alpha(n)$  is usually called the *inverse Ackerman function*.<sup>6</sup> Our earlier definition is equivalent to  $\alpha(n) = \min\{c \geq 1 \mid 2 \uparrow^{c+2} 3 \geq n\}$ ; in other words,  $\alpha(n) + 2$  is the inverse of the third column in the Ackermann hierarchy. The function  $\alpha(n)$  grows *much* more slowly than  $\log^{*c} n$  for any fixed  $c$ ; we have  $\alpha(n) \leq 3$  for all even *remotely imaginable* values of  $n$ . Nevertheless, the function  $\alpha(n)$  is eventually larger than any constant, so it is *not*  $O(1)$ .

$2 \uparrow^c n$	$n = 1$	$n = 2$	$n = 3$	$n = 4$	$n = 5$
$2n$	2	4	6	8	10
$2 \uparrow n$	2	4	8	16	32
$2 \uparrow\uparrow n$	2	4	16	65536	$2^{65536}$
$2 \uparrow\uparrow\uparrow n$	2	4	$65536$	$2^{2^{2^{\cdot^{\cdot^2}}}}_{65536}$	$2^{2^{2^{\cdot^{\cdot^2}}}}_{2^{2^{2^{\cdot^2}}}}_{65536}$
$2 \uparrow\uparrow\uparrow\uparrow n$	2	4	$2^{2^{2^{\cdot^{\cdot^2}}}}_{65536}$	$2^{2^{2^{\cdot^2}}}_{2^{2^{\cdot^2}}}_{65536}$	$2^{2^{2^{\cdot^2}}}_{2^{2^{2^{\cdot^2}}}}_{65536}$
$2 \uparrow\uparrow\uparrow\uparrow\uparrow n$	2	4	$2^{2^{2^{\cdot^2}}}_{65536}$	$2^{2^{2^{\cdot^2}}}_{2^{2^{\cdot^2}}}_{65536}$	$\langle\langle$ Yeah, right. $\rangle\rangle$
				$\langle\langle$ Very funny. $\rangle\rangle$	$\langle\langle$ Argh! My eyes! $\rangle\rangle$

Small (!! values of Ackermann's functions.

## 16.7 To infinity... and beyond!

Of course, one can generalize the inverse Ackermann function to functions that grow arbitrarily more slowly, starting with the *iterated* inverse Ackermann function

$$\alpha^*(n) = \begin{cases} 1 & \text{if } n \leq 4 \\ 1 + \alpha^*(\alpha(n)) & \text{otherwise,} \end{cases}$$

then the *iterated iterated* inverse Ackermann function

$$\alpha^{**}(n) = \begin{cases} 1 & \text{if } n \leq 4 \\ 1 + \alpha^{**}(\alpha(n)) & \text{otherwise,} \end{cases}$$

and then the diagonalized inverse Ackermann function

$$\text{Head-asplode}(n) = \min\{c \geq 1 \mid \alpha^{*c} n \leq 4\},$$

and so on forever. Fortunately(?), such functions appear extremely rarely in algorithm analysis. In fact, the only example (as far as Jeff knows) of a naturally-occurring super-constant sub-inverse-Ackermann function is a recent result of Seth Pettie<sup>7</sup>, who proved that if a splay tree is used as a double-ended queue — insertions and deletions of only smallest or largest elements — then the amortized cost of any operation is  $O(\alpha^*(n))$ .

<sup>6</sup>Strictly speaking, the name ‘inverse Ackerman function’ is inaccurate. One good formal definition of the true inverse Ackerman function is  $\tilde{\alpha}(n) = \min\{c \geq 1 \mid \lg^{*c} n \leq c\} = \min\{c \geq 1 \mid 2 \uparrow^{c+2} c \geq n\}$ . However, it’s not hard to prove that  $\tilde{\alpha}(n) \leq \alpha(n) \leq \tilde{\alpha}(n) + 1$  for all sufficiently large  $n$ , so the inaccuracy is completely forgivable. As I said in the previous footnote, the exact details of the definition are surprisingly irrelevant!

<sup>7</sup>Splay trees, Davenport-Schinzel sequences, and the deque conjecture. *Proceedings of the 19th Annual ACM-SIAM Symposium on Discrete Algorithms*, 1115–1124, 2008.

## Exercises

1. Consider the following solution for the union-find problem, called *union-by-weight*. Each set leader  $\bar{x}$  stores the number of elements of its set in the field  $weight(\bar{x})$ . Whenever we UNION two sets, the leader of the *smaller* set becomes a new child of the leader of the *larger* set (breaking ties arbitrarily).

<b>MAKESET(<math>x</math>):</b> $parent(x) \leftarrow x$ $weight(x) \leftarrow 1$	<b>UNION(<math>x, y</math>)</b> $\bar{x} \leftarrow \text{FIND}(x)$ $\bar{y} \leftarrow \text{FIND}(y)$ if $weight(\bar{x}) > weight(\bar{y})$ $parent(\bar{y}) \leftarrow \bar{x}$ $weight(\bar{x}) \leftarrow weight(\bar{x}) + weight(\bar{y})$ else $parent(\bar{x}) \leftarrow \bar{y}$ $weight(\bar{x}) \leftarrow weight(\bar{x}) + weight(\bar{y})$
<b>FIND(<math>x</math>):</b> while $x \neq parent(x)$ $x \leftarrow parent(x)$ return $x$	

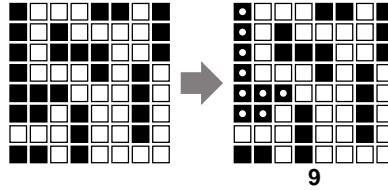
Prove that if we use union-by-weight, the *worst-case* running time of FIND( $x$ ) is  $O(\log n)$ , where  $n$  is the cardinality of the set containing  $x$ .

2. Consider a union-find data structure that uses union by depth (or equivalently union by rank) *without* path compression. For all integers  $m$  and  $n$  such that  $m \geq 2n$ , prove that there is a sequence of  $n$  MakeSet operations, followed by  $m$  UNION and FIND operations, that require  $\Omega(m \log n)$  time to execute.
3. Consider an arbitrary sequence of  $m$  MAKESET operations, followed by  $u$  UNION operations, followed by  $f$  FIND operations, and let  $n = m + u + f$ . Prove that if we use union by rank and FIND with path compression, all  $n$  operations are executed in  $O(n)$  time.
4. Describe and analyze a data structure to support the following operations on an array  $X[1..n]$  as quickly as possible. Initially,  $X[i] = 0$  for all  $i$ .
- Given an index  $i$  such that  $X[i] = 0$ , set  $X[i]$  to 1.
  - Given an index  $i$ , return  $X[i]$ .
  - Given an index  $i$ , return the smallest index  $j \geq i$  such that  $X[j] = 0$ , or report that no such index exists.

For full credit, the first two operations should run in *worst-case constant* time, and the amortized cost of the third operation should be as small as possible.

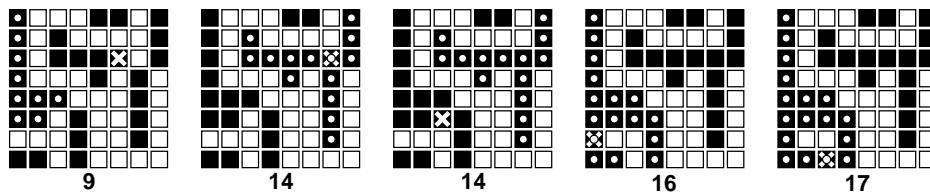
5. (a) Describe and analyze an algorithm to compute the size of the largest connected component of black pixels in an  $n \times n$  bitmap  $B[1..n, 1..n]$ .

For example, given the bitmap below as input, your algorithm should return the number 9, because the largest connected black component (marked with white dots on the right) contains nine pixels.



- (b) Design and analyze an algorithm  $\text{BLACKEN}(i, j)$  that colors the pixel  $B[i, j]$  black and returns the size of the largest black component in the bitmap. For full credit, the *amortized* running time of your algorithm (starting with an all-white bitmap) must be as small as possible.

For example, at each step in the sequence below, we blacken the pixel marked with an X. The largest black component is marked with white dots; the number underneath shows the correct output of the  $\text{BLACKEN}$  algorithm.



- (c) What is the *worst-case* running time of your  $\text{BLACKEN}$  algorithm?

6. Consider the following game. I choose a positive integer  $n$  and keep it secret; your goal is to discover this integer. We play the game in rounds. In each round, you write a list of *at most*  $n$  integers on the blackboard. If you write more than  $n$  numbers in a single round, you lose. (Thus, in the first round, you must write only the number 1; do you see why?) If  $n$  is one of the numbers you wrote, you win the game; otherwise, I announce which of the numbers you wrote is smaller or larger than  $n$ , and we proceed to the next round. For example:

You	Me
1	It's bigger than 1.
4, 42	It's between 4 and 42.
8, 15, 16, 23, 30	It's between 8 and 15.
9, 10, 11, 12, 13, 14	It's 11; you win!

Describe a strategy that allows you to win in  $O(\alpha(n))$  rounds!

*Obie looked at the seein' eye dog. Then at the twenty-seven 8 by 10 color glossy pictures with the circles and arrows and a paragraph on the back of each one... and then he looked at the seein' eye dog. And then at the twenty-seven 8 by 10 color glossy pictures with the circles and arrows and a paragraph on the back of each one and began to cry.*

*Because Obie came to the realization that it was a typical case of American blind justice, and there wasn't nothin' he could do about it, and the judge wasn't gonna look at the twenty-seven 8 by 10 color glossy pictures with the circles and arrows and a paragraph on the back of each one explainin' what each one was, to be used as evidence against us.*

*And we was fined fifty dollars and had to pick up the garbage. In the snow.*

*But that's not what I'm here to tell you about.*

— Arlo Guthrie, "Alice's Restaurant" (1966)

*I study my Bible as I gather apples.  
First I shake the whole tree, that the ripest might fall.  
Then I climb the tree and shake each limb,  
and then each branch and then each twig,  
and then I look under each leaf.*

— Martin Luther

## 17 Basic Graph Properties

### 17.1 Definitions

A **graph**  $G$  is a pair of sets  $(V, E)$ .  $V$  is a set of arbitrary objects that we call **vertices**<sup>1</sup> or **nodes**.  $E$  is a set of vertex pairs, which we call **edges** or occasionally **arcs**. In an *undirected* graph, the edges are unordered pairs, or just sets of two vertices; I will usually write  $uv$  instead of  $\{u,v\}$  to denote the undirected edge between  $u$  and  $v$ . In a *directed* graph, the edges are ordered pairs of vertices; I will usually write  $u \rightarrow v$  instead of  $(u,v)$  to denote the directed edge from  $u$  to  $v$ . We will usually be concerned only with **simple** graphs, where there is no edge from a vertex to itself and there is at most one edge from any vertex to any other.

Following standard (but admittedly confusing) practice, I'll also use  $V$  to denote the *number* of vertices in a graph, and  $E$  to denote the *number* of edges. Thus, in an undirected graph, we have  $0 \leq E \leq \binom{V}{2}$ , and in a directed graph,  $0 \leq E \leq V(V - 1)$ .

If  $(u, v)$  is an edge in an undirected graph, then  $u$  is a **neighbor** of  $v$  and vice versa. The **degree** of a node is the number of neighbors. In directed graphs, we have two kinds of neighbors. If  $u \rightarrow v$  is a directed edge, then  $u$  is a **predecessor** of  $v$  and  $v$  is a **successor** of  $u$ . The **in-degree** of a node is the number of predecessors, which is the same as the number of edges going into the node. The **out-degree** is the number of successors, or the number of edges going out of the node.

A graph  $G' = (V', E')$  is a **subgraph** of  $G = (V, E)$  if  $V' \subseteq V$  and  $E' \subseteq E$ .

A **path** is a sequence of edges, where each successive pair of edges shares a vertex, and all other edges are disjoint. A graph is **connected** if there is a path from any vertex to any other vertex. A disconnected graph consists of several **components**, which are its maximal connected subgraphs. Two vertices are in the same component if and only if there is a path between them.

---

<sup>1</sup>The singular of 'vertices' is **vertex**. The singular of 'matrices' is **matrix**. Unless you're speaking Italian, there is no such thing as a vertice, a matrice, an indice, an appendice, a helice, an apice, a vortice, a radice, a simplice, a codice, a directrice, a dominatrice, a Unice, a Kleenice, an Asterice, an Obelice, a Dogmatice, a Getafice, a Cacofonice, a Vitalstatistice, a Geriatrice, or Jimi Hendrice! You *will* lose points for using any of these so-called words.

A **cycle** is a path that starts and ends at the same vertex, and has at least one edge. An undirected graph is **acyclic** if no subgraph is a cycle; acyclic graphs are also called **forests**. **Trees** are special graphs that can be defined in several different ways. You can easily prove by induction (hint, hint, hint) that the following definitions are equivalent.

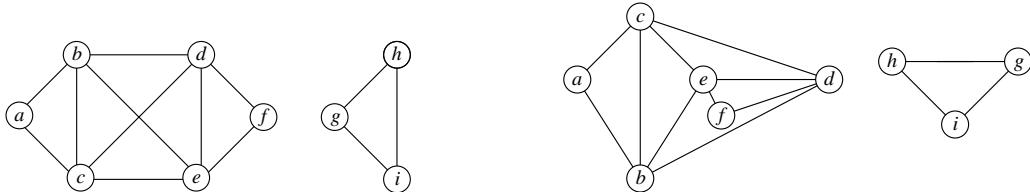
- A tree is a connected acyclic graph.
- A tree is a connected component of a forest.
- A tree is a connected graph with *at most*  $V - 1$  edges.
- A tree is a minimal connected graph; removing any edge makes the graph disconnected.
- A tree is an acyclic graph with *at least*  $V - 1$  edges.
- A tree is a maximal acyclic graph; adding an edge between any two vertices creates a cycle.

A **spanning tree** of a graph  $G$  is a subgraph that is a tree and contains every vertex of  $G$ . Of course, a graph can only have a spanning tree if it's connected. A **spanning forest** of  $G$  is a collection of spanning trees, one for each connected component of  $G$ .

Directed graphs can contain directed paths and directed cycles. A directed graph is **strongly connected** if there is a directed path from any vertex to any other. A directed graph is **acyclic** if it does not contain a directed cycle; directed acyclic graphs are often called **dags**.

## 17.2 Abstract Representations and Examples

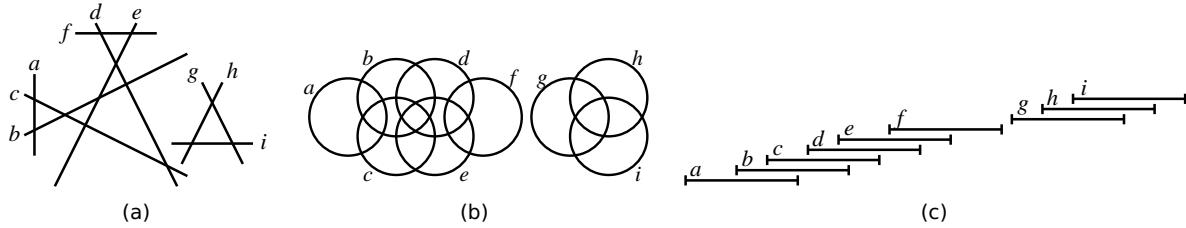
The most common way to visually represent graphs is by looking at an **embedding**. An embedding of a graph maps each vertex to a point in the plane and each edge to a curve or straight line segment between the two vertices. A graph is **planar** if it has an embedding where no two edges cross. The same graph can have many different embeddings, so it is important not to confuse a particular embedding with the graph itself. In particular, planar graphs can have non-planar embeddings!



A non-planar embedding of a planar graph with nine vertices, thirteen edges, and two connected components, and a planar embedding of the same graph.

However, embeddings are not the only useful representation of graphs. For example, the **intersection graph** of a collection of objects has a node for every object and an edge for every intersecting pair. Whether a particular graph can be represented as an intersection graph depends on what kind of object you want to use for the vertices. Different types of objects—line segments, rectangles, circles, etc.—define different classes of graphs. One particularly useful type of intersection graph is an **interval graph**, whose vertices are intervals on the real line, with an edge between any two intervals that overlap.

Another good example is the **dependency graph** of a recursive algorithm. Dependency graphs are directed acyclic graphs. The vertices are all the distinct recursive subproblems that arise when executing the algorithm on a particular input. There is an edge from one subproblem to another if evaluating



The example graph is also the intersection graph of (a) a set of line segments, (b) a set of circles, or (c) a set of intervals on the real line (stacked for visibility).

the second subproblem requires a recursive evaluation of the first subproblem. For example, for the Fibonacci recurrence

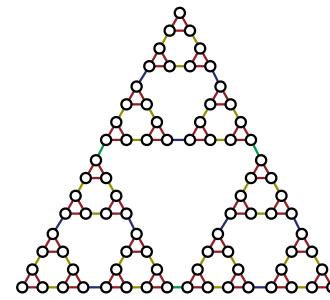
$$F_n = \begin{cases} 0 & \text{if } n = 0, \\ 1 & \text{if } n = 1, \\ F_{n-1} + F_{n-2} & \text{otherwise,} \end{cases}$$

the vertices of the dependency graph are the integers  $0, 1, 2, \dots, n$ , and for each integer  $i$  from 2 to  $n$ , and the edges are  $(i-1) \rightarrow i$  and  $(i-2) \rightarrow i$  for every integer  $i$  between 2 and  $n$ . For the edit distance recurrence

$$Edit(i, j) = \begin{cases} i & \text{if } j = 0 \\ j & \text{if } i = 0 \\ \min \left\{ \begin{array}{l} Edit(i-1, j) + 1, \\ Edit(i, j-1) + 1, \\ Edit(i-1, j-1) + [A[i] \neq B[j]] \end{array} \right\} & \text{otherwise} \end{cases}$$

the dependency graph is an  $m \times n$  grid with diagonals. Dynamic programming works efficiently for any recurrence that has a small dependency graph; a proper evaluation order ensures that each subproblem is visited *after* its predecessors.

A slightly more frivolous example of a graph is the **configuration graph** of a game, puzzle, or mechanism like tic-tac-toe, checkers, the Rubik's Cube, the Towers of Hanoi, or a Turing machine. The vertices of the configuration graph are all the valid configurations of the puzzle; there is an edge from one configuration to another if it is possible to transform one configuration into the other with a simple move. (Obviously, the precise definition depends on what moves are allowed.) Even for reasonably simple mechanisms, the configuration graph can be extremely complex, and we typically only have access to local information about the graph.



The configuration graph of the 4-disk Tower of Hanoi

Finally, the **finite-state automata** used in formal language theory are just labeled directed graphs. A deterministic finite-state automaton is usually formally defined as a 5-tuple  $M = (Q, \Sigma, \delta, q_0, A)$ , where  $Q$

is a finite set of *states*,  $\Sigma$  is a finite set called the *alphabet*,  $\delta : Q \times \Sigma \rightarrow Q$  is a *transition function*,  $q_0 \in Q$  is the *initial state*, and  $F \subseteq Q$  is the set of *accepting states*. But it is often useful to think of  $M$  as a directed graph  $G_M$  whose vertices are the states  $Q$ , and whose edges have the form  $q \xrightarrow{x} \delta(q, x)$  for every state  $q \in Q$  and character  $x \in \Sigma$ . Then basic questions about the language accepted by  $M$  can be phrased as questions about the graph  $G_M$ . For example, the language accepted by  $M$  is empty if and only if there is no path in  $G_M$  from  $q_0$  to an accepting state.

It's important not to confuse these examples/representations of graphs with the actual *definition*: A graph is a pair of sets  $(V, E)$ , where  $V$  is an arbitrary finite set, and  $E$  is a set of pairs (either ordered or unordered) of elements of  $V$ .

### 17.3 Graph Data Structures

There are two common data structures used to explicitly represent graphs: *adjacency matrices*<sup>2</sup> and *adjacency lists*.

The **adjacency matrix** of a graph  $G$  is a  $V \times V$  matrix, in which each entry indicates whether a particular edge is or is not in the graph:

$$A[i, j] := [(i, j) \in E].$$

For undirected graphs, the adjacency matrix is always *symmetric*:  $A[i, j] = A[j, i]$ . Since we don't allow edges from a vertex to itself, the diagonal elements  $A[i, i]$  are all zeros.

Given an adjacency matrix, we can decide in  $\Theta(1)$  time whether two vertices are connected by an edge just by looking in the appropriate slot in the matrix. We can also list all the neighbors of a vertex in  $\Theta(V)$  time by scanning the corresponding row (or column). This is optimal in the worst case, since a vertex can have up to  $V - 1$  neighbors; however, if a vertex has few neighbors, we may still have to examine every entry in the row to see them all. Similarly, adjacency matrices require  $\Theta(V^2)$  space, regardless of how many edges the graph actually has, so it is only space-efficient for very *dense* graphs.

	a	b	c	d	e	f	g	h	i
a	0	1	1	0	0	0	0	0	0
b	1	0	1	1	1	0	0	0	0
c	1	1	0	1	1	0	0	0	0
d	0	1	1	0	1	1	0	0	0
e	0	1	1	1	0	1	0	0	0
f	0	0	0	1	1	0	0	0	0
g	0	0	0	0	0	0	1	0	0
h	0	0	0	0	0	0	1	0	1
i	0	0	0	0	0	0	1	1	0

```

graph LR
    a --> b
    a --> c
    b --> a
    b --> c
    b --> d
    c --> a
    c --> b
    c --> d
    d --> b
    d --> c
    d --> e
    e --> b
    e --> c
    e --> d
    f --> d
    f --> e
    g --> h
    h --> g
    h --> i
    i --> g
    i --> h
  
```

Adjacency matrix and adjacency list representations for the example graph.

For *sparse* graphs—graphs with relatively few edges—adjacency lists are usually a better choice. An **adjacency list** is an array of linked lists, one list per vertex. Each linked list stores the neighbors of the corresponding vertex.

For undirected graphs, each edge  $(u, v)$  is stored twice, once in  $u$ 's neighbor list and once in  $v$ 's neighbor list; for directed graphs, each edge is stored only once. Either way, the overall space required for an adjacency list is  $O(V + E)$ . Listing the neighbors of a node  $v$  takes  $O(1 + \deg(v))$  time; just scan the neighbor list. Similarly, we can determine whether  $(u, v)$  is an edge in  $O(1 + \deg(u))$  time by scanning the neighbor list of  $u$ . For undirected graphs, we can speed up the search by simultaneously scanning

<sup>2</sup>See footnote 1.

the neighbor lists of both  $u$  and  $v$ , stopping either we locate the edge or when we fall off the end of a list. This takes  $O(1 + \min\{\deg(u), \deg(v)\})$  time.

The adjacency list structure should immediately remind you of hash tables with chaining. Just as with hash tables, we can make adjacency list structure more efficient by using something besides a linked list to store the neighbors. For example, if we use a hash table with constant load factor, when we can detect edges in  $O(1)$  expected time, just as with an adjacency list. In practice, this is only useful for vertices with very large degree, because the constant overhead in both the space and search time is larger for hash tables than for simple linked lists.

At this point, you might reasonably ask why anyone would ever use an adjacency matrix. After all, if you use hash tables to store the neighbors of each vertex, you can do everything as fast or faster with an adjacency list as with an adjacency matrix, only using less space. One compelling answer is that many graphs are *implicitly* represented by adjacency matrices. For example, intersection graphs are usually represented as an array or list of the underlying geometric objects. As long as we can test whether two objects overlap in constant time, we can apply any graph algorithm to an intersection graph by *pretending* that it is stored explicitly as an adjacency matrix. On the other hand, any data structure build from records with pointers between them can be seen as a directed graph. Algorithms for searching graphs can be applied to these data structures by *pretending* that the graph is represented explicitly using an adjacency list. Similarly, we can apply any graph algorithm to a configuration graph *as though* it were given to us as an adjacency list, provided we can enumerate all possible moves from a given configuration in constant time each.

To keep things simple, we'll consider only undirected graphs for the rest of this lecture, although the algorithms I'll describe also work for directed graphs.

## 17.4 Traversing connected graphs

Suppose we want to visit every node in a connected graph (represented either explicitly or implicitly). The simplest method to do this is an algorithm called *depth-first search*, which can be written either recursively or iteratively. It's exactly the same algorithm either way; the only difference is that we can actually see the 'recursion' stack in the non-recursive version. Both versions are initially passed a *source* vertex  $s$ .

```
RECURSIVEDFS( $v$ ):
    if  $v$  is unmarked
        mark  $v$ 
        for each edge  $vw$ 
            RECURSIVEDFS( $w$ )
```

```
ITERATIVEDFS( $s$ ):
    PUSH( $s$ )
    while the stack is not empty
         $v \leftarrow \text{POP}$ 
        if  $v$  is unmarked
            mark  $v$ 
            for each edge  $vw$ 
                PUSH( $w$ )
```

Depth-first search is one (perhaps the most common) instance of a general family of graph traversal algorithms. The generic graph traversal algorithm stores a set of candidate edges in some data structure that I'll call a 'bag'. The only important properties of a 'bag' are that we can put stuff into it and then later take stuff back out. (In C++ terms, think of the 'bag' as a template for a real data structure.) A stack is a particular type of bag, but certainly not the only one. Here is the generic traversal algorithm:

```
TRAVERSE( $s$ ):
    put  $s$  in bag
    while the bag is not empty
        take  $v$  from the bag
        if  $v$  is unmarked
            mark  $v$ 
            for each edge  $vw$ 
                put  $w$  into the bag
```

This traversal algorithm clearly marks each vertex in the graph *at most* once. In order to show that it visits every node in the graph *at least* once, we modify the algorithm slightly; the modifications are highlighted in red. Instead of keeping vertices in the bag, the modified algorithm stores pairs of vertices. This allows us to remember, whenever we visit a vertex  $v$  for the first time, which previously-visited vertex  $p$  put  $v$  into the bag. This vertex  $p$  is called the *parent* of  $v$ .

```
TRAVERSE( $s$ ):
    put  $(\emptyset, s)$  in bag
    while the bag is not empty
        take  $(p, v)$  from the bag (*) 
        if  $v$  is unmarked
            mark  $v$ 
             $\text{parent}(v) \leftarrow p$ 
            for each edge  $vw$  (†)
                put  $(v, w)$  into the bag (**)
```

**Lemma 1.** *TRAVERSE( $s$ ) marks every vertex in any connected graph exactly once, and the set of pairs  $(v, \text{parent}(v))$  with  $\text{parent}(v) \neq \emptyset$  defines a spanning tree of the graph.*

**Proof:** The algorithm obviously marks  $s$ . Let  $v$  be any vertex other than  $s$ , and let  $s \rightarrow \dots \rightarrow u \rightarrow v$  be the directed path from  $s$  to  $v$  with the minimum number of edges. Since the graph is connected, such a path always exists. (If  $s$  and  $v$  are neighbors, then  $u = s$ , and the path has just one edge.) If the algorithm marks  $u$ , then it must put  $(u, v)$  into the bag, so it must later take  $(u, v)$  out of the bag, at which point  $v$  must be marked. Thus, by induction on the shortest-path distance from  $s$ , the algorithm marks every vertex in the graph. This implies that  $\text{parent}(v)$  is well-defined for every vertex  $v$ .

The algorithm clearly marks every vertex at most once, so it must mark every vertex *exactly* once.

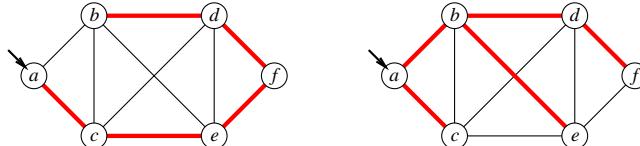
Call any pair  $(v, \text{parent}(v))$  with  $\text{parent}(v) \neq \emptyset$  a *parent edge*. For any node  $v$ , the directed path of parent edges  $v \rightarrow \text{parent}(v) \rightarrow \text{parent}(\text{parent}(v)) \rightarrow \dots$  eventually leads back to  $s$ , so the set of parent edges form a connected graph. Clearly, both endpoints of every parent edge are marked, and the number of parent edges is exactly one less than the number of vertices. Thus, the parent edges form a spanning tree.  $\square$

The exact running time of the traversal algorithm depends on how the graph is represented and what data structure is used as the ‘bag’, but we can make a few general observations. Because each vertex is marked at most once, the for loop (†) is executed at most  $V$  times. Each edge  $uv$  is put into the bag exactly twice; once as the pair  $(u, v)$  and once as the pair  $(v, u)$ , so line (\*\*) is executed at most  $2E$  times. Finally, we can’t take more things out of the bag than we put in, so line (\*) is executed at most  $2E + 1$  times.

## 17.5 Examples

Let’s first assume that the graph is represented by an adjacency list, so that the overhead of the for loop (†) is only constant time per edge.

- If we implement the ‘bag’ using a *stack*, we recover our original depth-first search algorithm. Each execution of  $(\star)$  or  $(\star\star)$  takes constant time, so the overall running time is  $O(V + E)$ . If the graph is connected, we have  $V \leq E + 1$ , and so we can simplify the running time to  $O(E)$ . The spanning tree formed by the parent edges is called a ***depth-first spanning tree***. The exact shape of the tree depends on the start vertex and on the order that neighbors are visited in the for loop  $(\dagger)$ , but in general, depth-first spanning trees are long and skinny.
- If we use a *queue* instead of a stack, we get ***breadth-first search***. Again, each execution of  $(\star)$  or  $(\star\star)$  takes constant time, so the overall running time for connected graphs is still  $O(E)$ . In this case, the ***breadth-first spanning tree*** formed by the parent edges contains *shortest paths* from the start vertex  $s$  to every other vertex in its connected component. We’ll see shortest paths again in a future lecture. Again, exact shape of a breadth-first spanning tree depends on the start vertex and on the order that neighbors are visited in the for loop  $(\dagger)$ , but in general, shortest path trees are short and bushy.



A depth-first spanning tree and a breadth-first spanning tree of one component of the example graph, with start vertex  $a$ .

- Now suppose the edges of the graph are weighted. If we implement the ‘bag’ using a *priority queue*, always extracting the minimum-weight edge in line  $(\star)$ , the resulting algorithm is reasonably called ***shortest-first search***. In this case, each execution of  $(\star)$  or  $(\star\star)$  takes  $O(\log E)$  time, so the overall running time is  $O(V + E \log E)$ , which simplifies to  $O(E \log E)$  if the graph is connected. For this algorithm, the set of parent edges form the ***minimum spanning tree*** of the connected component of  $s$ . Surprisingly, as long as all the edge weights are distinct, the resulting tree does *not* depend on the start vertex or the order that neighbors are visited; in this case, there is only one minimum spanning tree. We’ll see minimum spanning trees again in the next lecture.

If the graph is represented using an adjacency matrix instead of an adjacency list, finding all the neighbors of each vertex in line  $(\dagger)$  takes  $O(V)$  time. Thus, depth- and breadth-first search each run in  $O(V^2)$  time, and ‘shortest-first search’ runs in  $O(V^2 + E \log E) = O(V^2 \log V)$  time.

## 17.6 Searching disconnected graphs

If the graph is disconnected, then  $\text{TRAVERSE}(s)$  only visits the nodes in the connected component of the start vertex  $s$ . If we want to visit all the nodes in every component, we can use the following ‘wrapper’ around our generic traversal algorithm. Since  $\text{TRAVERSE}$  computes a spanning tree of one component,  $\text{TRAVERSEALL}$  computes a spanning forest of the entire graph.

$\text{TRAVERSEALL}(s):$ for all vertices $v$ if $v$ is unmarked $\text{TRAVERSE}(v)$
--

Some textbooks claim that this wrapper can only be used with depth-first search; they’re wrong.

## Exercises

1. Prove that the following definitions are all equivalent.
  - A tree is a connected acyclic graph.
  - A tree is a connected component of a forest.
  - A tree is a connected graph with *at most*  $V - 1$  edges.
  - A tree is a minimal connected graph; removing any edge makes the graph disconnected.
  - A tree is an acyclic graph with *at least*  $V - 1$  edges.
  - A tree is a maximal acyclic graph; adding an edge between any two vertices creates a cycle.
  
2. Prove that any connected acyclic graph with  $n \geq 2$  vertices has at least two vertices with degree 1. Do not use the words ‘tree’ or ‘leaf’, or any well-known properties of trees; your proof should follow entirely from the definitions.
  
3. Let  $G$  be a connected graph, and let  $T$  be a depth-first spanning tree of  $G$  rooted at some node  $v$ . Prove that if  $T$  is also a breadth-first spanning tree of  $G$  rooted at  $v$ , then  $G = T$ .
  
4. Whenever groups of pigeons gather, they instinctively establish a *pecking order*. For any pair of pigeons, one pigeon always pecks the other, driving it away from food or potential mates. The same pair of pigeons always chooses the same pecking order, even after years of separation, no matter what other pigeons are around. Surprisingly, the overall pecking order can contain cycles—for example, pigeon  $A$  pecks pigeon  $B$ , which pecks pigeon  $C$ , which pecks pigeon  $A$ .
  - (a) Prove that any finite set of pigeons can be arranged in a row from left to right so that every pigeon pecks the pigeon immediately to its left. Pretty please.
  - (b) Suppose you are given a directed graph representing the pecking relationships among a set of  $n$  pigeons. The graph contains one vertex per pigeon, and it contains an edge  $i \rightarrow j$  if and only if pigeon  $i$  pecks pigeon  $j$ . Describe and analyze an algorithm to compute a pecking order for the pigeons, as guaranteed by part (a).
  
5. You are helping a group of ethnographers analyze some oral history data they have collected by interviewing members of a village to learn about the lives of people lived there over the last two hundred years. From the interviews, you have learned about a set of people, all now deceased, whom we will denote  $P_1, P_2, \dots, P_n$ . The ethnographers have collected several facts about the lifespans of these people. Specifically, for some pairs  $(P_i, P_j)$ , the ethnographers have learned one of the following facts:
  - (a)  $P_i$  died before  $P_j$  was born.
  - (b)  $P_i$  and  $P_j$  were both alive at some moment.

Naturally, the ethnographers are not sure that their facts are correct; memories are not so good, and all this information was passed down by word of mouth. So they’d like you to determine whether the data they have collected is at least internally consistent, in the sense that there could have existed a set of people for which all the facts they have learned simultaneously hold.

Describe and analyze an algorithm to answer the ethnographers' problem. Your algorithm should either output possible dates of birth and death that are consistent with all the stated facts, or it should report correctly that no such dates exist.

6. Let  $G = (V, E)$  be a given directed graph.
  - (a) The *transitive closure*  $G^T$  is a directed graph with the same vertices as  $G$ , that contains any edge  $u \rightarrow v$  if and only if there is a directed path from  $u$  to  $v$  in  $G$ . Describe an efficient algorithm to compute the transitive closure of  $G$ .
  - (b) The *transitive reduction*  $G^{TR}$  is the smallest graph (meaning fewest edges) whose transitive closure is  $G^T$ . Describe an efficient algorithm to compute the transitive reduction of  $G$ .
7. A graph  $(V, E)$  is *bipartite* if the vertices  $V$  can be partitioned into two subsets  $L$  and  $R$ , such that every edge has one vertex in  $L$  and the other in  $R$ .
  - (a) Prove that every tree is a bipartite graph.
  - (b) Describe and analyze an efficient algorithm that determines whether a given undirected graph is bipartite.
8. An *Euler tour* of a graph  $G$  is a closed walk through  $G$  that traverses every edge of  $G$  exactly once.
  - (a) Prove that a connected graph  $G$  has an Euler tour if and only if every vertex has even degree.
  - (b) Describe and analyze an algorithm to compute an Euler tour in a given graph, or correctly report that no such graph exists.
9. The  $d$ -dimensional hypercube is the graph defined as follows. There are  $2d$  vertices, each labeled with a different string of  $d$  bits. Two vertices are joined by an edge if their labels differ in exactly one bit.
  - (a) A Hamiltonian cycle in a graph  $G$  is a cycle of edges in  $G$  that visits every vertex of  $G$  exactly once. Prove that for all  $d \geq 2$ , the  $d$ -dimensional hypercube has a Hamiltonian cycle.
  - (b) Which hypercubes have an Euler tour (a closed walk that traverses every edge exactly once)?  
*[Hint: This is very easy.]*
10. **Racetrack** (also known as *Graph Racers* and *Vector Rally*) is a two-player paper-and-pencil racing game that Jeff played on the bus in 5th grade.<sup>3</sup> The game is played with a track drawn on a sheet of graph paper. The players alternately choose a sequence of grid points that represent the motion of a car around the track, subject to certain constraints explained below.

Each car has a *position* and a *velocity*, both with integer  $x$ - and  $y$ -coordinates. The initial position is a point on the starting line, chosen by the player; the initial velocity is always  $(0, 0)$ . At each step, the player optionally increments or decrements either or both coordinates of the car's

---

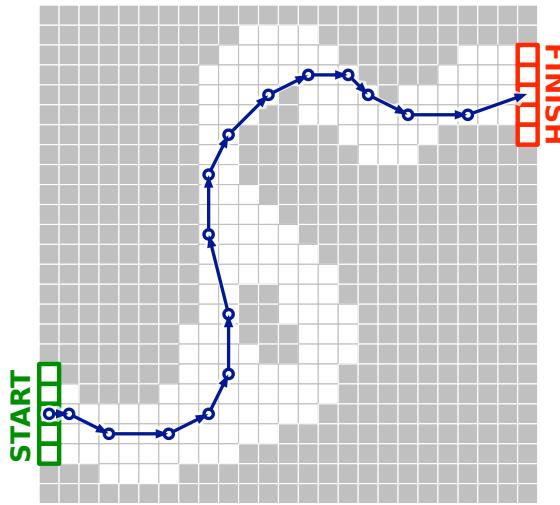
<sup>3</sup>The actual game is a bit more complicated than the version described here. In particular, in the actual game, the boundaries of the track are a free-form curve, and (at least by default) the entire line segment between any two consecutive positions must lie inside the track. In the version Jeff played, if a car does run off the track, the car starts its next turn with zero velocity, at the legal grid point closest to where the car left the track.

velocity; in other words, each component of the velocity can change by at most 1 in a single step. The car's new position is then determined by adding the new velocity to the car's previous position. The new position must be inside the track; otherwise, the car crashes and that player loses the race. The race ends when the first car reaches a position *on* the finish line.

Suppose the racetrack is represented by an  $n \times n$  array of bits, where each 0 bit represents a grid point inside the track, each 1 bit represents a grid point outside the track, the 'starting line' is the first column, and the 'finish line' is the last column.

Describe and analyze an algorithm to find the minimum number of steps required to move a car from the starting line to the finish line of a given racetrack. [Hint: Build a graph. What are the vertices? What are the edges? What problem is this?]

velocity	position
(0, 0)	(1, 5)
(1, 0)	(2, 5)
(2, -1)	(4, 4)
(3, 0)	(7, 4)
(2, 1)	(9, 5)
(1, 2)	(10, 7)
(0, 3)	(10, 10)
(-1, 4)	(9, 14)
(0, 3)	(9, 17)
(1, 2)	(10, 19)
(2, 2)	(12, 21)
(2, 1)	(14, 22)
(2, 0)	(16, 22)
(1, -1)	(17, 21)
(2, -1)	(19, 20)
(3, 0)	(22, 20)
(3, 1)	(25, 21)



A 16-step Racetrack run, on a  $25 \times 25$  track. This is *not* the shortest run on this track.

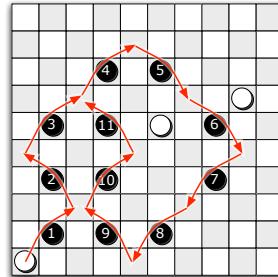
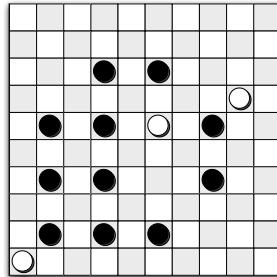
- \*11. Draughts/checkers is a game played on an  $m \times m$  grid of squares, alternately colored light and dark. (The game is usually played on an  $8 \times 8$  or  $10 \times 10$  board, but the rules easily generalize to any board size.) Each dark square is occupied by at most one game piece (usually called a *checker* in the U.S.), which is either black or white; light squares are always empty. One player ('White') moves the white pieces; the other ('Black') moves the black pieces.

Consider the following simple version of the game, essentially American checkers or British draughts, but where every piece is a king.<sup>4</sup> Pieces can be moved in any of the four diagonal directions, either one or two steps at a time. On each turn, a player either *moves* one of her pieces one step diagonally into an empty square, or makes a series of *jumps* with one of her checkers. In a single jump, a piece moves to an empty square two steps away in any diagonal direction, but only if the intermediate square is occupied by a piece of the opposite color; this enemy piece is *captured* and immediately removed from the board. Multiple jumps are allowed in a single turn as long as they are made by the same piece. A player wins if her opponent has no pieces left on the board.

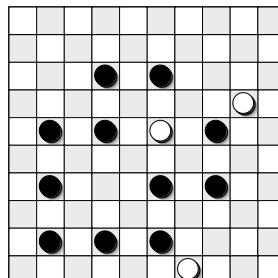
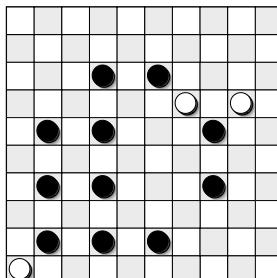
Describe an algorithm that correctly determines whether White can capture every black piece, thereby winning the game, *in a single turn*. The input consists of the width of the board ( $m$ ), a list of positions of white pieces, and a list of positions of black pieces. For full credit, your algorithm

<sup>4</sup>Most other variants of draughts have 'flying kings', which behave *very* differently than what's described here.

should run in  $O(n)$  time, where  $n$  is the total number of pieces. [Hint: The greedy strategy—make arbitrary jumps until you get stuck—does **not** always find a winning sequence of jumps even when one exists. See problem 8. Parity, parity, parity.]



White wins in one turn.



White cannot win in one turn from either of these positions.

*We must all hang together, gentlemen, or else we shall most assuredly hang separately.*

— Benjamin Franklin, at the signing of the Declaration of Independence (July 4, 1776)

*It is a very sad thing that nowadays there is so little useless information.*

— Oscar Wilde

*A ship in port is safe, but that is not what ships are for.*

— Rear Admiral Grace Murray Hopper

## 18 Minimum Spanning Trees

### 18.1 Introduction

Suppose we are given a connected, undirected, *weighted* graph. This is a graph  $G = (V, E)$  together with a function  $w: E \rightarrow \mathbb{R}$  that assigns a *weight*  $w(e)$  to each edge  $e$ . For this lecture, we'll assume that the weights are real numbers. Our task is to find the **minimum spanning tree** of  $G$ , that is, the spanning tree  $T$  minimizing the function

$$w(T) = \sum_{e \in T} w(e).$$

To keep things simple, I'll assume that all the edge weights are distinct:  $w(e) \neq w(e')$  for any pair of edges  $e$  and  $e'$ . Distinct weights guarantee that the minimum spanning tree of the graph is unique. Without this condition, there may be several different minimum spanning trees. For example, if all the edges have weight 1, then *every* spanning tree is a minimum spanning tree with weight  $V - 1$ .

If we have an algorithm that assumes the edge weights are unique, we can still use it on graphs where multiple edges have the same weight, as long as we have a consistent method for breaking ties. One way to break ties consistently is to use the following algorithm in place of a simple comparison. **SHORTEREDGE** takes as input four integers  $i, j, k, l$ , and decides which of the two edges  $(i, j)$  and  $(k, l)$  has ‘smaller’ weight.

```

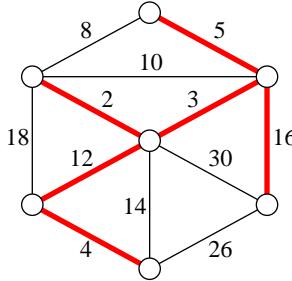
SHORTEREDGE( $i, j, k, l$ )
  if  $w(i, j) < w(k, l)$  return  $(i, j)$ 
  if  $w(i, j) > w(k, l)$  return  $(k, l)$ 
  if  $\min(i, j) < \min(k, l)$  return  $(i, j)$ 
  if  $\min(i, j) > \min(k, l)$  return  $(k, l)$ 
  if  $\max(i, j) < \max(k, l)$  return  $(i, j)$ 
  {{if  $\max(i, j) < \max(k, l)$  }} return  $(k, l)$ 

```

### 18.2 The Only Minimum Spanning Tree Algorithm

There are several different methods for computing minimum spanning trees, but really they are all instances of the following generic algorithm. The situation is similar to the previous lecture, where we saw that depth-first search and breadth-first search were both instances of a single generic traversal algorithm.

The generic minimum spanning tree algorithm maintains an acyclic subgraph  $F$  of the input graph  $G$ , which we will call an *intermediate spanning forest*.  $F$  is a subgraph of the minimum spanning tree of  $G$ , and every component of  $F$  is a minimum spanning tree of its vertices. Initially,  $F$  consists of  $n$  one-node trees. The generic algorithm merges trees together by adding certain edges between them. When the



A weighted graph and its minimum spanning tree.

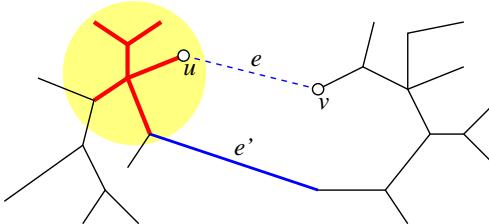
algorithm halts,  $F$  consists of a single  $n$ -node tree, which must be the minimum spanning tree. Obviously, we have to be careful about *which* edges we add to the evolving forest, since not every edge is in the minimum spanning tree.

The intermediate spanning forest  $F$  induces two special types of edges. An edge is *useless* if it is not an edge of  $F$ , but both its endpoints are in the same component of  $F$ . For each component of  $F$ , we associate a *safe* edge—the minimum-weight edge with exactly one endpoint in that component. Different components might or might not have different safe edges. Some edges are neither safe nor useless—we call these edges *undecided*.

All minimum spanning tree algorithms are based on two simple observations.<sup>1</sup>

**Lemma 1.** *The minimum spanning tree contains every safe edge.*

**Proof:** Let  $T$  be the minimum spanning tree. Suppose  $F$  has a ‘bad’ component whose safe edge  $e = (u, v)$  is not in  $T$ . Since  $T$  is connected, it contains a unique path from  $u$  to  $v$ , and at least one edge  $e'$  on this path has exactly one endpoint in the bad component. Removing  $e'$  from the minimum spanning tree and adding  $e$  gives us a new spanning tree. Since  $e$  is the bad component’s safe edge, we have  $w(e') > w(e)$ , so the new spanning tree has smaller total weight than  $T$ . But this is impossible— $T$  is the *minimum* spanning tree. So  $T$  must contain every safe edge.  $\square$



Proving that every safe edge is in the minimum spanning tree. The ‘bad’ component of  $F$  is highlighted.

**Lemma 2.** *The minimum spanning tree contains no useless edge.*

**Proof:** Adding any useless edge to  $F$  would introduce a cycle.  $\square$

Our generic minimum spanning tree algorithm repeatedly adds one or more safe edges to the evolving forest  $F$ . Whenever we add new edges to  $F$ , some undecided edges become safe, and others become useless. To specify a particular algorithm, we must decide which safe edges to add, and we must describe how to identify new safe and new useless edges, at each iteration of our generic template.

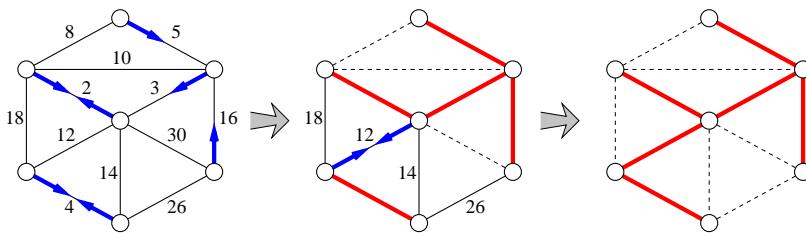
<sup>1</sup>Lemmas 1 and 2 are actually special cases of two more general observations. First, for any partition of the vertices of  $G$  into two disjoint subsets, the minimum-weight edge with one endpoint in each subset is in the minimum spanning tree. Second, the maximum-weight edge in any cycle in  $G$  is **not** in the minimum spanning tree.

### 18.3 Boruvka's Algorithm

The oldest and arguably simplest minimum spanning tree algorithm was discovered by Boruvka in 1926, long before computers even existed, and practically before the invention of graph theory!<sup>2</sup> The algorithm was rediscovered by Choquet in 1938; again by Florek, Łukaziewicz, Perkal, Stienhaus, and Zubrzycki in 1951; and again by Sollin some time in the early 1960s. Because Sollin was the only Western computer scientist in this list—Choquet was a civil engineer; Florek and his co-authors were anthropologists—this is often called ‘Sollin’s algorithm’, especially in the parallel computing literature.

The Boruvka/Choquet/Florek/Łukaziewicz/Perkal/Stienhaus/Zubrzycki/Sollin algorithm can be summarized in one line:

BORUVKA: Add all the safe edges and recurse.



Boruvka's algorithm run on the example graph. Thick edges are in  $F$ . Arrows point along each component's safe edge. Dashed edges are useless.

We can find all the safe edge in the graph in  $O(E)$  time as follows. First, we count the components of  $F$  using whatever-first search, using the TRAVERSEALL wrapper function. As we count, we label every vertex with its component number; that is, every vertex in the first component traversed by TRAVERSEALL gets label 1, every vertex in the second component gets label 2, and so on.

If  $F$  has only one component, we're done. Otherwise, we compute an array  $S[1..V]$  of edges, where  $S[i]$  is the minimum-weight edge with one endpoint in the  $i$ th component (or a sentinel value NULL if there are less than  $i$  components). To compute this array, we consider each edge  $uv$  in the input graph  $G$ . If the endpoints  $u$  and  $v$  have the same label, then  $uv$  is useless. Otherwise, we compare the weight of  $uv$  to the weights of  $S[\text{label}(u)]$  and  $S[\text{label}(v)]$  and update the array entries if necessary.

```
BORUVKA( $V, E$ ):
   $F = (V, \emptyset)$ 
  TRAVERSEALL( $F$ )  «count and label components»
  while  $F$  has more than one component
    ADDALLSAFEEDGES( $V, E$ )
    TRAVERSEALL( $F$ )
  return  $F$ 
```

```
ADDALLSAFEEDGES( $V, E$ ):
  for  $i \leftarrow 1$  to  $V$ 
     $S[i] \leftarrow \text{NULL}$  «sentinel: } $w(\text{NULL}) := \infty$ »
  for each edge  $uv \in E$ 
    if  $\text{label}(u) \neq \text{label}(v)$ 
      if  $w(uv) < w(S[\text{label}(u)])$ 
         $S[\text{label}(u)] \leftarrow uv$ 
      if  $w(uv) < w(S[\text{label}(v)])$ 
         $S[\text{label}(v)] \leftarrow uv$ 
  for  $i \leftarrow 1$  to  $V$ 
    if  $S[i] \neq \text{NULL}$ 
      add  $S[i]$  to  $F$ 
```

Each call to TRAVERSEALL requires  $O(V)$  time, because the forest  $F$  has at most  $V - 1$  edges. Assuming the graph is represented by an adjacency list, the rest of each iteration of the main while loop requires

<sup>2</sup>Leonard Euler published the first graph theory result, his famous theorem about the bridges of Königsburg, in 1736. However, the first textbook on graph theory, written by Dénes König, was not published until 1936.

$O(E)$  time, because we spend constant time on each edge. Because the graph is connected, we have  $V \leq E + 1$ , so each iteration of the while loop takes  $O(E)$  time.

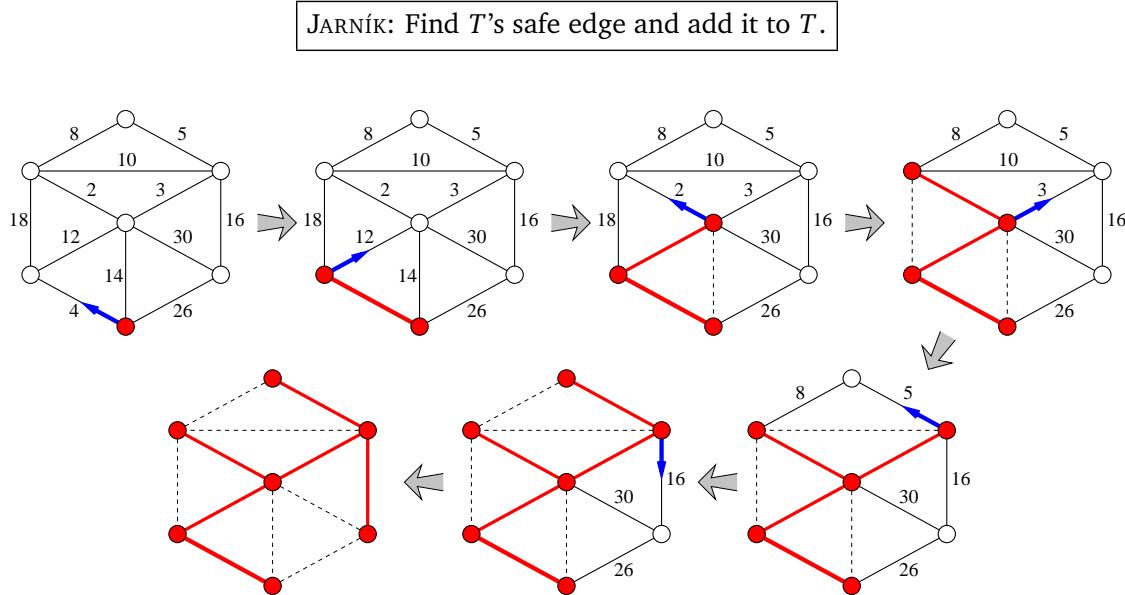
Each iteration reduces the number of components of  $F$  by at least a factor of two—the worst case occurs when the components coalesce in pairs. Since  $F$  initially has  $V$  components, the while loop iterates at most  $O(\log V)$  times. Thus, the overall running time of Boruvka's algorithm is  $O(E \log V)$ .

Despite its relatively obscure origin, early algorithms researchers were aware of Boruvka's algorithm, but dismissed it as being ‘too complicated’! As a result, despite its simplicity and efficiency, Boruvka's algorithm is rarely mentioned in algorithms and data structures textbooks. On the other hand, more recent algorithms to compute minimum spanning trees are all generalizations of Boruvka's algorithm, not the other two classical algorithms described next.

#### 18.4 Jarník's ('Prim's') Algorithm

The next oldest minimum spanning tree algorithm was first described by the Czech mathematician Vojtěch Jarník in a 1929 letter to Boruvka; Jarník published his discovery the following year. The algorithm was independently rediscovered by Kruskal in 1956, by Prim in 1957, by Loberman and Weinberger in 1957, and finally by Dijkstra in 1958. Prim, Loberman, Weinberger, and Dijkstra all (eventually) knew of and even cited Kruskal's paper, but since Kruskal also described two other minimum-spanning-tree algorithms in the same paper, *this* algorithm is usually called ‘Prim's algorithm’, or sometimes ‘the Prim/Dijkstra algorithm’, even though by 1958 Dijkstra already had another algorithm (inappropriately) named after him.

In Jarník's algorithm, the forest  $F$  contains only one nontrivial component  $T$ ; all the other components are isolated vertices. Initially,  $T$  consists of an arbitrary vertex of the graph. The algorithm repeats the following step until  $T$  spans the whole graph:



Jarník's algorithm run on the example graph, starting with the bottom vertex.  
At each stage, thick edges are in  $T$ , an arrow points along  $T$ 's safe edge, and dashed edges are useless.

To implement Jarník's algorithm, we keep all the edges adjacent to  $T$  in a priority queue. When we pull the minimum-weight edge out of the priority queue, we first check whether both of its endpoints are in  $T$ . If not, we add the edge to  $T$  and then add the new neighboring edges to the priority queue. In other words, Jarník's algorithm is another instance of the generic graph traversal algorithm we saw last

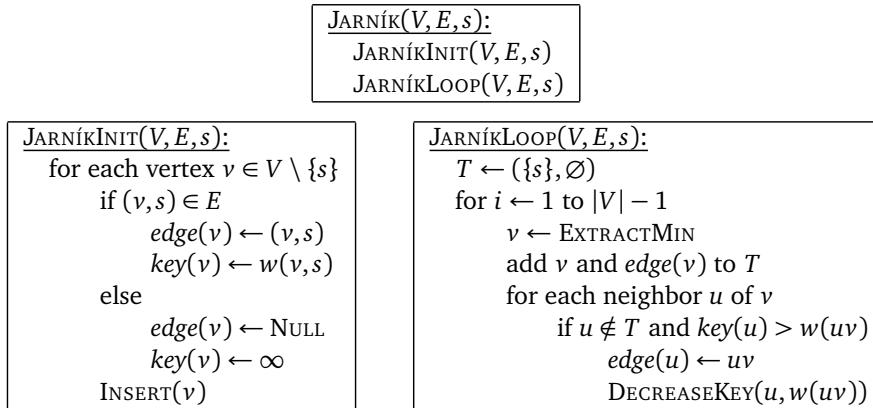
time, using a priority queue as the ‘bag’! If we implement the algorithm this way, the algorithm runs in  $O(E \log E) = O(E \log V)$  time.

### \*18.5 Improving Jarník’s Algorithm

We can improve Jarník’s algorithm using a more advanced priority queue data structure called a *Fibonacci heap*, first described by Michael Fredman and Robert Tarjan in 1984. Fibonacci heaps support the standard priority queue operations `INSERT`, `EXTRACTMIN`, and `DECREASEKEY`. However, unlike standard binary heaps, which require  $O(\log n)$  time for every operation, Fibonacci heaps support `INSERT` and `DECREASEKEY` in constant *amortized* time. The amortized cost of `EXTRACTMIN` is still  $O(\log n)$ .

To apply this faster data structure, we keep *vertices* in the priority queue instead of edge, where the key for each vertex  $v$  is either the minimum-weight edge between  $v$  and the evolving tree  $T$ , or  $\infty$  if there is no such edge. We can `INSERT` all the vertices into the priority queue at the beginning of the algorithm; then, whenever we add a new edge to  $T$ , we may need to decrease the keys of some neighboring vertices.

To make the description easier, we break the algorithm into two parts. `JARNÍKINIT` initializes the priority queue; `JARNÍKLOOP` is the main algorithm. The input consists of the vertices and edges of the graph, plus the start vertex  $s$ . For each vertex  $v$ , we maintain both its key  $key(v)$  and the incident edge  $edge(v)$  such that  $w(edge(v)) = key(v)$ .



The operations `INSERT` and `EXTRACTMIN` are each called  $O(V)$  times once for each vertex except  $s$ , and `DECREASEKEY` is called  $O(E)$  times, at most twice for each edge. Thus, if we use a Fibonacci heap, the improved algorithm runs in  $O(E + V \log V)$  time, which is faster than Boruvka’s algorithm unless  $E = O(V)$ .

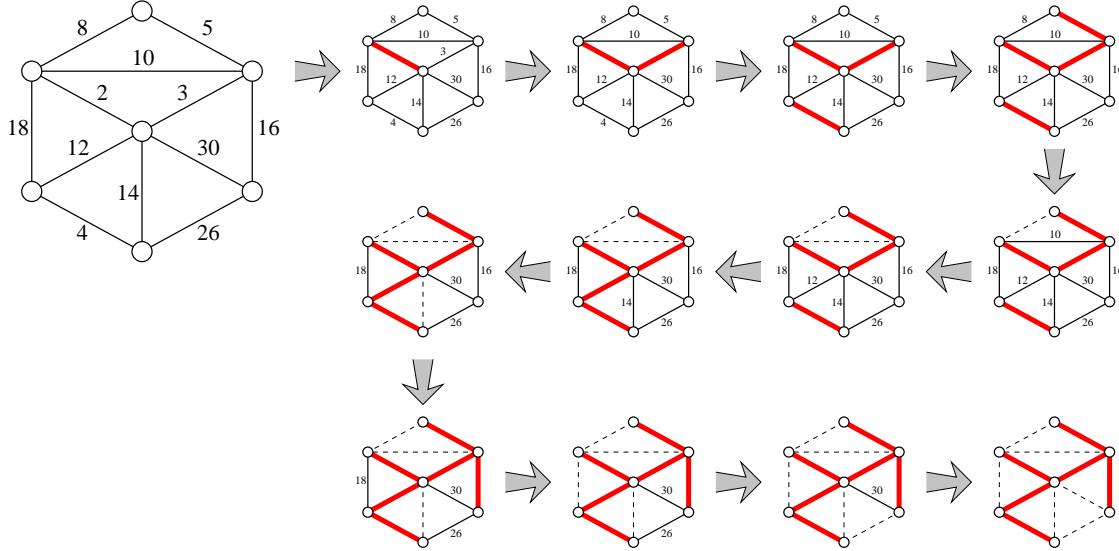
In practice, however, this improvement is rarely faster than the naive implementation using a binary heap, unless the graph is extremely large and dense. The Fibonacci heap algorithms are quite complex, and the hidden constants in both the running time and space are significant—not outrageous, but certainly bigger than the hidden constant 1 in the  $O(\log n)$  time bound for binary heap operations.

### 18.6 Kruskal’s Algorithm

The last minimum spanning tree algorithm I’ll discuss was first described by Kruskal in 1956, in the same paper where he rediscovered Jarník’s algorithm. Kruskal was motivated by ‘a typewritten translation (of obscure origin)’ of Boruvka’s original paper, claiming that Boruvka’s algorithm was ‘unnecessarily

elaborate.<sup>3</sup> This algorithm was also rediscovered in 1957 by Loberman and Weinberger, but somehow avoided being renamed after them.

KRUSKAL: Scan all edges in increasing weight order; if an edge is safe, add it to  $F$ .



Kruskal's algorithm run on the example graph. Thick edges are in  $F$ . Dashed edges are useless.

Since we examine the edges in order from lightest to heaviest, any edge we examine is safe if and only if its endpoints are in different components of the forest  $F$ . To prove this, suppose the edge  $e$  joins two components  $A$  and  $B$  but is not safe. Then there would be a lighter edge  $e'$  with exactly one endpoint in  $A$ . But this is impossible, because (inductively) any previously examined edge has both endpoints in the same component of  $F$ .

Just as in Boruvka's algorithm, each component of  $F$  has a 'leader' node. An edge joins two components of  $F$  if and only if the two endpoints have different leaders. But unlike Boruvka's algorithm, we do not recompute leaders from scratch every time we add an edge. Instead, when two components are joined, the two leaders duke it out in a nationally-televised no-holds-barred steel-cage grudge match.<sup>4</sup> One of the two emerges victorious as the leader of the new larger component. More formally, we will use our earlier algorithms for the UNION-FIND problem, where the vertices are the elements and the components of  $F$  are the sets. Here's a more formal description of the algorithm:

**KRUSKAL( $V, E$ ):**

- sort  $E$  by increasing weight
- $F \leftarrow (V, \emptyset)$
- for each vertex  $v \in V$ 
  - $\text{MAKESET}(v)$
- for  $i \leftarrow 1$  to  $|E|$ 
  - $uv \leftarrow i$ th lightest edge in  $E$
  - if  $\text{FIND}(u) \neq \text{FIND}(v)$ 
    - $\text{UNION}(u, v)$
    - add  $uv$  to  $F$
- return  $F$

<sup>3</sup>To be fair, Boruvka's original paper was unnecessarily elaborate, but in his followup paper, also published in 1927, simplified his algorithm to its current modern form. Kruskal was apparently unaware of Boruvka's second paper. Stupid Iron Curtain.

<sup>4</sup>Live at the Assembly Hall! Only \$49.95 on Pay-Per-View!

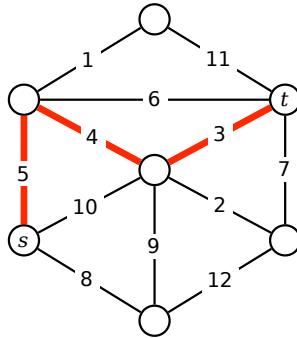
In our case, the sets are components of  $F$ , and  $n = V$ . Kruskal's algorithm performs  $O(E)$  FIND operations, two for each edge in the graph, and  $O(V)$  UNION operations, one for each edge in the minimum spanning tree. Using union-by-rank and path compression allows us to perform each UNION or FIND in  $O(\alpha(E, V))$  time, where  $\alpha$  is the not-quite-constant inverse-Ackerman function. So ignoring the cost of sorting the edges, the running time of this algorithm is  $O(E \alpha(E, V))$ .

We need  $O(E \log E) = O(E \log V)$  additional time just to sort the edges. Since this is bigger than the time for the UNION-FIND data structure, the overall running time of Kruskal's algorithm is  $O(E \log V)$ , exactly the same as Boruvka's algorithm, or Jarník's algorithm with a normal (non-Fibonacci) heap.

## Exercises

1. Most classical minimum-spanning-tree algorithms use the notions of ‘safe’ and ‘useless’ edges described in the lecture notes, but there is an alternate formulation. Let  $G$  be a weighted undirected graph, where the edge weights are distinct. We say that an edge  $e$  is *dangerous* if it is the longest edge in some cycle in  $G$ , and *useful* if it does not lie in any cycle in  $G$ .
  - (a) Prove that the minimum spanning tree of  $G$  contains every useful edge.
  - (b) Prove that the minimum spanning tree of  $G$  does not contain any dangerous edge.
  - (c) Describe and analyze an efficient implementation of the “anti-Kruskal” MST algorithm: Examine the edges of  $G$  in *decreasing* order; if an edge is dangerous, remove it from  $G$ . [Hint: It won’t be as fast as Kruskal’s algorithm.]
2. Let  $G = (V, E)$  be an arbitrary connected graph with weighted edges.
  - (a) Prove that for any partition of the vertices  $V$  into two subsets, the minimum-weight edge with one endpoint in each subset is in the minimum spanning tree of  $G$ .
  - (b) Prove that the maximum-weight edge in any cycle of  $G$  is *not* in the minimum spanning tree of  $G$ .
  - (c) Prove or disprove: The minimum spanning tree of  $G$  includes the minimum-weighted edge in *every* cycle in  $G$ .
3. Throughout this lecture note, we assumed that no two edges in the input graph have equal weights, which implies that the minimum spanning tree is unique. In fact, a weaker condition on the edge weights implies MST uniqueness.
  - (a) Describe an edge-weighted graph that has a unique minimum spanning tree, even though two edges have equal weights.
  - (b) Prove that an edge-weighted graph  $G$  has a *unique* minimum spanning tree if and only if the following conditions hold:
    - For any partition of the vertices of  $G$  into two subsets, the minimum-weight edge with one endpoint in each subset is unique.
    - The maximum-weight edge in any cycle of  $G$  is unique.
  - (c) Describe and analyze an algorithm to determine whether or not a graph has a unique minimum spanning tree.

4. (a) Describe and analyze an algorithm to compute the *maximum*-weight spanning tree of a given edge-weighted graph.
- (b) A *feedback edge set* of a graph  $G$  is a subset  $F$  of the edges such that every cycle in  $G$  contains at least one edge in  $F$ . In other words, removing every edge in  $F$  makes the graph  $G$  acyclic. Describe and analyze a fast algorithm to compute the minimum weight feedback edge set of a given edge-weighted graph.
5. Consider a path between two vertices  $s$  and  $t$  in an undirected weighted graph  $G$ . The *bottleneck length* of this path is the maximum weight of any edge in the path. The *bottleneck distance* between  $s$  and  $t$  is the minimum bottleneck length of any path from  $s$  to  $t$ . (If there are no paths from  $s$  to  $t$ , the bottleneck distance between  $s$  and  $t$  is  $\infty$ .) Describe and analyze an algorithm to



The bottleneck distance between  $s$  and  $t$  is 5.

compute the bottleneck distance between *every* pair of vertices in an arbitrary undirected weighted graph. Assume that no two edges have the same weight.

6. Suppose you are given a graph  $G$  with weighted edges and a minimum spanning tree  $T$  of  $G$ .
- (a) Describe an algorithm to update the minimum spanning tree when the weight of a single edge  $e$  is decreased.
- (b) Describe an algorithm to update the minimum spanning tree when the weight of a single edge  $e$  is increased.

In both cases, the input to your algorithm is the edge  $e$  and its new weight; your algorithms should modify  $T$  so that it is still a minimum spanning tree. [Hint: Consider the cases  $e \in T$  and  $e \notin T$  separately.]

7. (a) Describe and analyze an algorithm to find the *second smallest spanning tree* of a given graph  $G$ , that is, the spanning tree of  $G$  with smallest total weight except for the minimum spanning tree.
- \*(b) Describe and analyze an efficient algorithm to compute, given a weighted undirected graph  $G$  and an integer  $k$ , the  $k$  spanning trees of  $G$  with smallest weight.
8. We say that a graph  $G = (V, E)$  is *dense* if  $|E| = \Theta(|V|^2)$ . Describe a modification of Jarník's minimum-spanning tree algorithm that runs in  $O(|V|^2)$  time (independent of  $E$ ) when the input graph is dense, using only simple data structures (and in particular, *without* using a Fibonacci heap).

9. Consider an algorithm that first performs  $k$  passes of Borùvka's algorithm, and then runs Jarník's algorithm (with a Fibonacci heap) on the resulting contracted graph.
  - (a) What is the running time of this hybrid algorithm, as a function of  $V$ ,  $E$ , and  $k$ ?
  - (b) For which value of  $k$  is this running time minimized? What is the resulting running time?
10. Describe an algorithm to compute the minimum spanning tree of an  $n$ -vertex *planar* graph in  $O(n)$  time. [*Hint: Contracting an edge in a planar graph yields another planar graph. Any planar graph with  $n$  vertices has at most  $3n - 6$  edges.*]

*Well, ya turn left by the fire station in the village and take the old post road by the reservoir and... no, that won't do.*

*Best to continue straight on by the tar road until you reach the schoolhouse and then turn left on the road to Bennett's Lake until... no, that won't work either.*

*East Millinocket, ya say? Come to think of it, you can't get there from here.*

— Robert Bryan and Marshall Dodge,  
*Bert and I and Other Stories from Down East* (1961)

*Hey farmer! Where does this road go?*

*Been livin' here all my life, it ain't gone nowhere yet.*

*Hey farmer! How do you get to Little Rock?*

*Listen stranger, you can't get there from here.*

*Hey farmer! You don't know very much do you?*

*No, but I ain't lost.*

— Michelle Shocked, "Arkansas Traveler" (1992)

## 19 Shortest Paths

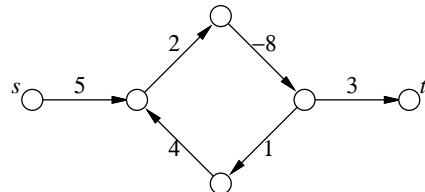
### 19.1 Introduction

Suppose we are given a weighted *directed* graph  $G = (V, E, w)$  with two special vertices, and we want to find the shortest path from a *source* vertex  $s$  to a *target* vertex  $t$ . That is, we want to find the directed path  $p$  starting at  $s$  and ending at  $t$  that minimizes the function

$$w(p) := \sum_{u \rightarrow v \in p} w(u \rightarrow v).$$

For example, if I want to answer the question 'What's the fastest way to drive from my old apartment in Champaign, Illinois to my wife's old apartment in Columbus, Ohio?', I might use a graph whose vertices are cities, edges are roads, weights are driving times,  $s$  is Champaign, and  $t$  is Columbus.<sup>1</sup> The graph is directed, because driving times along the same road might be different in different directions.<sup>2</sup>

Perhaps counter to intuition, we will allow the weights on the edges to be negative. Negative edges make our lives complicated, because the presence of a negative cycle might imply that there is no shortest path. In general, a shortest path from  $s$  to  $t$  exists if and only if there is *at least one* path from  $s$  to  $t$ , but there is no path from  $s$  to  $t$  that touches a negative cycle. If there is a negative cycle between  $s$  and  $t$ , then we can always find a shorter path by going around the cycle one more time.



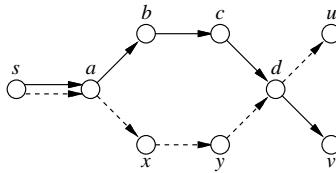
There is no shortest path from  $s$  to  $t$ .

<sup>1</sup>West on Church, north on Prospect, east on I-74, south on I-465, east on Airport Expressway, north on I-65, east on I-70, north on Grandview, east on 5th, north on Olentangy River, east on Dodridge, north on High, west on Kelso, south on Neil. Depending on traffic. We both live in Urbana now.

<sup>2</sup>At one time, there was a speed trap on I-70 just east of the Indiana/Ohio border, but only for eastbound traffic.

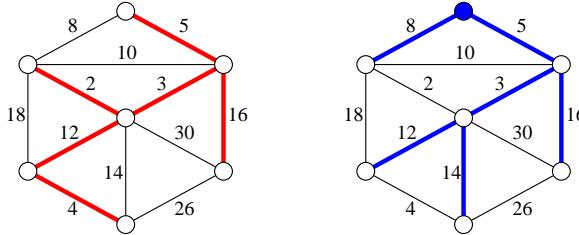
Almost every algorithm known for solving this problem actually solves (large portions of) the following more general *single source shortest path* or *SSSP* problem: find the shortest path from the source vertex  $s$  to every other vertex in the graph. In fact, the problem is usually solved by finding a *shortest path tree* rooted at  $s$  that contains all the desired shortest paths.

It's not hard to see that if shortest paths are unique, then they form a tree. To prove this, it's enough to observe that any subpath of a shortest path is also a shortest path. If there are multiple shortest paths to some vertices, we can always choose one shortest path to each vertex so that the union of the paths is a tree. If there are shortest paths to two vertices  $u$  and  $v$  that diverge, then meet, then diverge again, we can modify one of the paths so that the two paths only diverge once.



If  $s \rightarrow a \rightarrow b \rightarrow c \rightarrow d \rightarrow v$  and  $s \rightarrow a \rightarrow x \rightarrow y \rightarrow d \rightarrow u$  are shortest paths,  
then  $s \rightarrow a \rightarrow b \rightarrow c \rightarrow d \rightarrow u$  is also a shortest path.

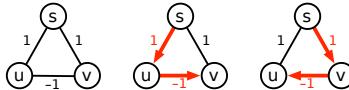
Shortest path trees and minimum spanning trees are very different. For one thing, there is only one minimum spanning tree (if edge weights are distinct), but each source vertex induces a different shortest path tree. Moreover, in general, *all* of these shortest path trees are different from the minimum spanning tree.



A minimum spanning tree and a shortest path tree (rooted at the topmost vertex) of the same graph.

## 19.2 Warning!

Throughout this lecture, we will explicitly consider *only* directed graphs. All of the algorithms described in this lecture also work for undirected graphs with some minor modifications, *but only if negative edges are prohibited*. Dealing with negative edges in undirected graphs is considerably more subtle. We cannot simply replace every undirected edge with a pair of directed edges, because this would transform any negative edge into a short negative cycle. Subpaths of an undirected shortest path that contains a negative edge are *not* necessarily shortest paths; consequently, the set of all undirected shortest paths from a single source vertex may not define a tree, even if shortest paths are unique.



An undirected graph where shortest paths from  $s$  are unique but do not define a tree.

A complete treatment of undirected graphs with negative edges is beyond the scope of this lecture (if not the entire course). I will only mention that a *single* shortest path in an undirected graph with negative edges can be computed in  $O(VE + V^2 \log V)$  time, by a reduction to maximum weighted matching.

### 19.3 The Only SSSP Algorithm

Just like graph traversal and minimum spanning trees, there are several different SSSP algorithms, but they are all special cases of a single generic algorithm, first proposed by Lester Ford in 1956, and independently by George Dantzig in 1957.<sup>3</sup> Each vertex  $v$  in the graph stores two values, which (inductively) describe a *tentative* shortest path from  $s$  to  $v$ .

- $dist(v)$  is the length of the tentative shortest  $s \leadsto v$  path, or  $\infty$  if there is no such path.
- $pred(v)$  is the predecessor of  $v$  in the tentative shortest  $s \leadsto v$  path, or `NULL` if there is no such vertex.

In fact, the predecessor pointers automatically define a tentative shortest path *tree*; they play the same role as the ‘parent’ pointers in our generic graph traversal algorithm.

At the beginning of the algorithm, we already know that  $dist(s) = 0$  and  $pred(s) = \text{NULL}$ . For every vertex  $v \neq s$ , we initially set  $dist(v) = \infty$  and  $pred(v) = \text{NULL}$  to indicate that we do not know of *any* path from  $s$  to  $v$ .

We call an edge  $u \rightarrow v$  *tense* if  $dist(u) + w(u \rightarrow v) < dist(v)$ . If  $u \rightarrow v$  is tense, then the tentative shortest path  $s \leadsto v$  is incorrect, since the path  $s \leadsto u \rightarrow v$  is shorter. Our generic algorithm repeatedly finds a tense edge in the graph and *relaxes* it:

```
RELAX( $u \rightarrow v$ ):
     $dist(v) \leftarrow dist(u) + w(u \rightarrow v)$ 
     $pred(v) \leftarrow u$ 
```

If there are no tense edges, our algorithm is finished, and we have our desired shortest path tree.

The correctness of the relaxation algorithm follows directly from four simple claims:

1. For every vertex  $v$ ,  $dist(v)$  is always greater than or equal to the shortest-path distance from  $s$  to  $v$ . This claim can be verified by induction on the number of relaxations.
2. If no edge is tense, then for every vertex  $v$ ,  $dist(v)$  is the shortest-path distance from  $s$  to  $v$ . This claim can be verified by induction on the length of the true shortest path, or on the rank of the true shortest path distance if all edge weights are non-negative.
3. If no edge is tense, then for every vertex  $v$ , the path  $s \rightarrow \dots \rightarrow pred(pred(v)) \rightarrow pred(v) \rightarrow v$  is a shortest path from  $s$  to  $v$ . This claim can be proved by induction on the number of edges in the path. (Why does  $pred(v)$  have its current value?)
4. The algorithm halts if and only if there is no negative cycle reachable from  $s$ . The ‘only if’ direction is easy—if there is a reachable negative cycle, then after the first edge in the cycle is relaxed, the cycle *always* has at least one tense edge. The ‘if’ direction follows from the fact that every relaxation step reduces either the number of vertices with  $dist(v) = \infty$  by 1 or reduces the sum of the finite shortest path lengths by the difference between two edge weights.

I haven’t said anything about how we detect which edges can be relaxed, or in what order we relax them. In order to make this easier, we can refine the relaxation algorithm slightly, into something closely resembling the generic graph traversal algorithm. We maintain a ‘bag’ of vertices, initially containing just the source vertex  $s$ . Whenever we take a vertex  $u$  out of the bag, we scan all of its outgoing edges, looking for something to relax. Whenever we successfully relax an edge  $u \rightarrow v$ , we put  $v$  into the bag. Unlike our generic graph traversal algorithm, the same vertex might be visited many times.

---

<sup>3</sup>Specifically, Dantzig showed that the shortest path problem can be phrased as a linear programming problem, and then described an interpretation of his simplex method in terms of the original graph. His description is equivalent to Ford’s relaxation strategy.

```
INITSSSP( $s$ ):
 $dist(s) \leftarrow 0$ 
 $pred(s) \leftarrow \text{NULL}$ 
for all vertices  $v \neq s$ 
 $dist(v) \leftarrow \infty$ 
 $pred(v) \leftarrow \text{NULL}$ 
```

```
GENERICSSSP( $s$ ):
INITSSSP( $s$ )
put  $s$  in the bag
while the bag is not empty
take  $u$  from the bag
for all edges  $u \rightarrow v$ 
if  $u \rightarrow v$  is tense
RELAX( $u \rightarrow v$ )
put  $v$  in the bag
```

Just as with graph traversal, using different data structures for the ‘bag’ gives us different algorithms. There are three obvious choices to try: a stack, a FIFO queue, and a priority queue. Unfortunately, if we use a stack, we have to perform  $\Theta(2^V)$  relaxation steps in the worst case! (Proving this is a good homework problem.) The other two possibilities prove to be much more efficient.

## 19.4 Dijkstra’s Algorithm

If we implement the bag using a priority queue, where the key of a vertex  $v$  is  $dist(v)$ , we obtain an algorithm first ‘published’<sup>4</sup> by Leyzorek, Gray, Johnson, Ladew, Meaker, Petry, and Seitz in 1957, and then later independently rediscovered by Edsger Dijkstra in 1959. A very similar algorithm was also described by George Dantzig in 1958.

Dijkstra’s algorithm, as it is universally known<sup>5</sup>, is particularly well-behaved if the graph has no negative-weight edges. In this case, it’s not hard to show (by induction, of course) that the vertices are scanned in increasing order of their shortest-path distance from  $s$ . It follows that each vertex is scanned at most once, and thus that each edge is relaxed at most once. Since the key of each vertex in the heap is its tentative distance from  $s$ , the algorithm performs a DECREASEKEY operation every time an edge is relaxed. Thus, the algorithm performs at most  $E$  DECREASEKEYS. Similarly, there are at most  $V$  INSERT and EXTRACTMIN operations. Thus, if we store the vertices in a Fibonacci heap, the total running time of Dijkstra’s algorithm is  $O(E + V \log V)$ ; if we use a regular binary heap, the running time is  $O(E \log V)$ .

This analysis assumes that no edge has negative weight. Dijkstra’s algorithm (in the form I’m presenting here) is still *correct* if there are negative edges<sup>6</sup>, but the worst-case running time could be exponential. (Proving this unfortunate fact is a good homework problem.) On the other hand, in practice, Dijkstra’s algorithm is usually quite fast even for graphs with negative edges.

## 19.5 The $A^*$ Heuristic

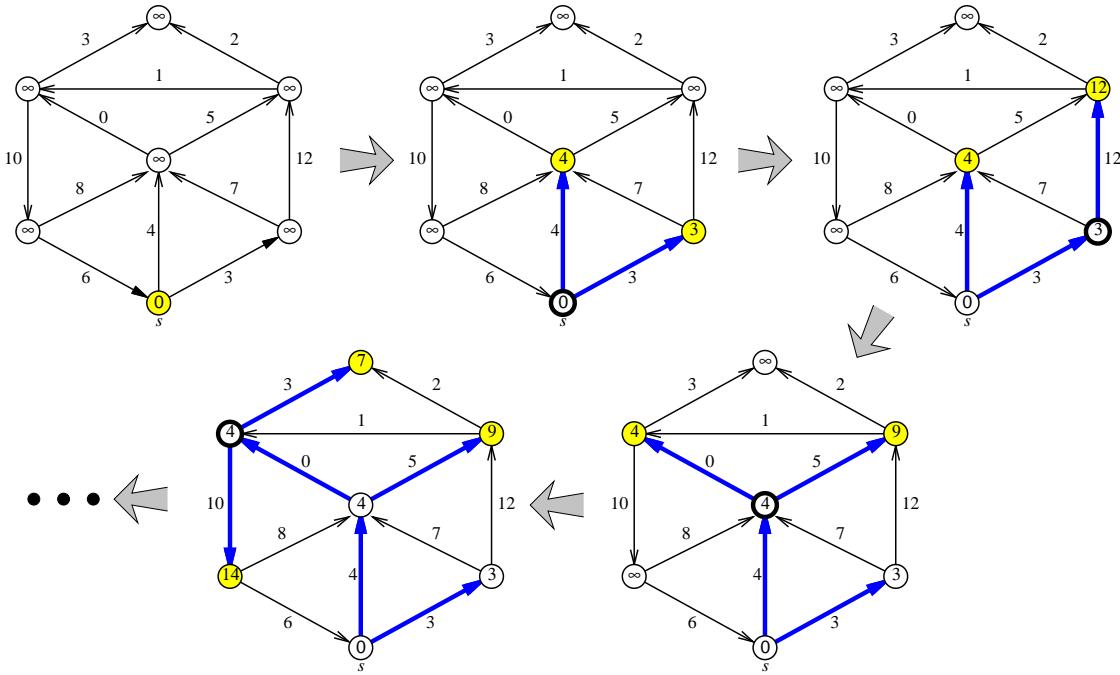
A slight generalization of Dijkstra’s algorithm, commonly known as the  $A^*$  algorithm, is frequently used to find a shortest path from a single source node  $s$  to a single target node  $t$ .  $A^*$  uses a black-box function GUESSEDISTANCE( $v, t$ ) that returns an estimate of the distance from  $v$  to  $t$ . The only difference between Dijkstra and  $A^*$  is that the key of a vertex  $v$  is  $dist(v) + \text{GUESSEDISTANCE}(v, t)$ .

The function GUESSEDISTANCE is called *admissible* if GUESSEDISTANCE( $v, t$ ) never overestimates the actual shortest path distance from  $v$  to  $t$ . If GUESSEDISTANCE is admissible and the actual edge weights are all non-negative, the  $A^*$  algorithm computes the actual shortest path from  $s$  to  $t$  at least as quickly as

<sup>4</sup>in the first annual report on a research project performed for the Combat Development Department of the Army Electronic Proving Ground

<sup>5</sup>I will follow this common convention, despite the historical inaccuracy, because I don’t think anybody wants to read about the “Leyzorek-Gray-Johnson-Ladew-Meaker-Petry-Seitz algorithm”.

<sup>6</sup>Most textbooks, and indeed Dijkstra’s original paper, present a version of Dijkstra’s algorithm that gives incorrect results for graphs with negative edges, because it *never* visits the same vertex more than once.



Four phases of Dijkstra's algorithm run on a graph with no negative edges.

At each phase, the shaded vertices are in the heap, and the bold vertex has just been scanned.

The bold edges describe the evolving shortest path tree.

Dijkstra's algorithm. The closer  $\text{GUESSESDISTANCE}(v, t)$  is to the real distance from  $v$  to  $t$ , the faster the algorithm. However, in the worst case, the running time is still  $O(E + V \log V)$ .

The heuristic is especially useful in situations where the actual graph is not known. For example,  $A^*$  can be used to solve many puzzles (15-puzzle, Freecell, Shanghai, Sokoban, Atomix, Rush Hour, Rubik's Cube, Racetrack, ...) and other path planning problems where the starting and goal configurations are given, but the graph of all possible configurations and their connections is not given explicitly.

## 19.6 Shimbel's Algorithm ('Bellman-Ford')

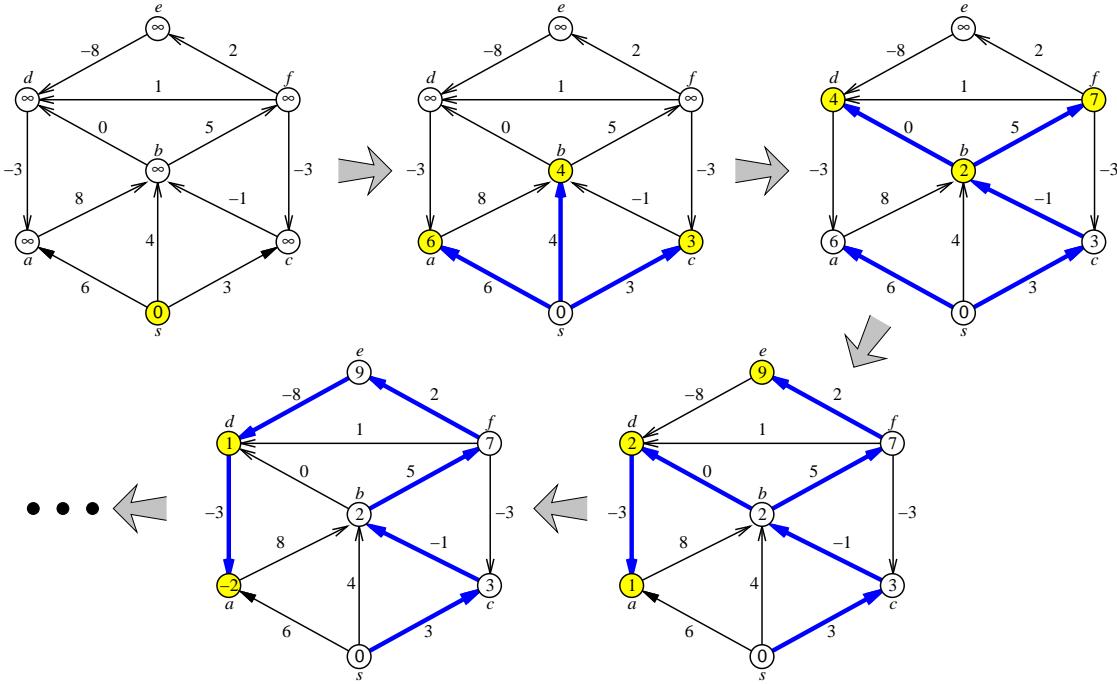
If we replace the heap in Dijkstra's algorithm with a FIFO queue, we obtain an algorithm first described by Shimbel in 1955, then independently rediscovered by Moore in 1957, by Woodbury and Dantzig in 1957, and by Bellman in 1958. Because Bellman explicitly used Ford's formulation of relaxing edges, this algorithm is usually called 'Bellman-Ford', although some early sources refer to 'Bellman-Shimbel'. Shimbel's algorithm is efficient even if there are negative edges, and it can be used to quickly detect the presence of negative cycles. If there are no negative edges, however, Dijkstra's algorithm is faster. (In fact, in practice, Dijkstra's algorithm is often faster even for graphs with negative edges.)

The easiest way to analyze the algorithm is to break the execution into *phases*, by introducing an imaginary *token*. Before we even begin, we insert the token into the queue. The current phase ends when we take the token out of the queue; we begin the next phase by reinserting the token into the queue. The 0th phase consists entirely of scanning the source vertex  $s$ . The algorithm ends when the queue contains *only* the token. A simple inductive argument (hint, hint) implies the following invariant:

At the end of the  $i$ th phase, for every vertex  $v$ ,  $\text{dist}(v)$  is at most the length of the shortest path  $s \rightsquigarrow v$  consisting of at most  $i$  edges.

Since a shortest path can only pass through each vertex once, either the algorithm halts before the

$V$ th phase, or the graph contains a negative cycle. In each phase, we scan each vertex at most once, so we relax each edge at most once, so the running time of a single phase is  $O(E)$ . Thus, the overall running time of Shimbel's algorithm is  $O(VE)$ .



Four phases of Shimbel's algorithm run on a directed graph with negative edges.

Nodes are taken from the queue in the order  $s \diamond a \ b \ c \ \diamond d \ f \ b \diamond a \ e \ d \diamond d \ a \ \diamond \diamond$ , where  $\diamond$  is the end-of-phase token. Shaded vertices are in the queue at the end of each phase. The bold edges describe the evolving shortest path tree.

Once we understand how the phases of Shimbel's algorithm behave, we can simplify the algorithm considerably. Instead of using a queue to perform a partial breadth-first search of the graph in each phase, we can simply scan through the adjacency list directly and try to relax every edge in the graph.

<b>SHIMBELSSSP(<math>s</math>)</b> INITSSSP( $s$ ) repeat $V$ times: for every edge $u \rightarrow v$ if $u \rightarrow v$ is tense RELAX( $u \rightarrow v$ ) for every edge $u \rightarrow v$ if $u \rightarrow v$ is tense return 'Negative cycle!'
--

This is how most textbooks present the ‘Bellman-Ford’ algorithm.<sup>7</sup> The  $O(VE)$  running time of this version of the algorithm should be obvious, but it may not be clear that the algorithm is still correct. To prove correctness, we have to show that our earlier invariant holds; as before, this can be proved by induction on  $i$ .

<sup>7</sup>In fact, this is closer to the description that Shimbel and Bellman used. Bob Tarjan recognized in the early 1980s that Shimbel's algorithm is equivalent to Dijkstra's algorithm with a queue instead of a heap.

## 19.7 Shimbel's Algorithm as Dynamic Programming

Shimbel's algorithm can also be recast as a dynamic programming algorithm. Let  $dist_i(v)$  denote the length of the shortest path  $s \rightsquigarrow v$  consisting of at most  $i$  edges. It's not hard to see that this function obeys the following recurrence:

$$dist_i(v) = \begin{cases} 0 & \text{if } i = 0 \text{ and } v = s \\ \infty & \text{if } i = 0 \text{ and } v \neq s \\ \min \left\{ \min_{u \rightarrow v \in E} (dist_{i-1}(u) + w(u \rightarrow v)) \right\} & \text{otherwise} \end{cases}$$

For the moment, let's assume the graph has no negative cycles; our goal is to compute  $dist_{V-1}(t)$ . We can clearly memoize this two-parameter function into a two-dimensional array. A straightforward dynamic programming evaluation of this recurrence looks like this:

```
SHIMBELDP(s)
  dist[0, s] ← 0
  for every vertex v ≠ s
    dist[0, v] ← ∞
  for i ← 1 to V - 1
    for every vertex v
      dist[i, v] ← dist[i - 1, v]
      for every edge u → v
        if dist[i, v] > dist[i - 1, u] + w(u → v)
          dist[i, v] ← dist[i - 1, u] + w(u → v)
```

Now let us make two minor changes to this algorithm. First, we remove one level of indentation from the last three lines. This may change the order in which we examine edges, but the modified algorithm still computes  $dist_i(v)$  for all  $i$  and  $v$ . Second, we change the indices in the last two lines from  $i - 1$  to  $i$ . This change may cause the distances  $dist[i, v]$  to approach the true shortest-path distances more quickly than before, but the algorithm is still correct.

```
SHIMBELDP2(s)
  dist[0, s] ← 0
  for every vertex v ≠ s
    dist[0, v] ← ∞
  for i ← 1 to V - 1
    for every vertex v
      dist[i, v] ← dist[i - 1, v]
      for every edge u → v
        if dist[i, v] > dist[i, u] + w(u → v)
          dist[i, v] ← dist[i, u] + w(u → v)
```

Now notice that the iteration index  $i$  is completely redundant! We really only need to keep a one-dimensional array of distances, which means we don't need to scan the vertices in each iteration of the main loop.

```
SHIMBELDP3(s)
  dist[s] ← 0
  for every vertex v ≠ s
    dist[v] ← ∞
  for i ← 1 to V - 1
    for every edge u → v
      if dist[v] > dist[u] + w(u → v)
        dist[v] ← dist[u] + w(u → v)
```

The resulting algorithm is (almost) identical to our earlier algorithm SHIMBEL! The first three lines initialize the shortest path distances, and the last two lines check whether an edge is tense, and if so, relaxes it. The only thing missing is the explicit maintenance of predecessors, but that's easy to add.

## 19.8 Greedy Relaxation?

Here's another algorithm that fits our generic framework, but which I've never seen analyzed.

Repeatedly relax the tensest edge.

Specifically, let's define the 'tension' of an edge  $u \rightarrow v$  as follows:

$$\text{tension}(u \rightarrow v) = \max \{0, \text{dist}(v) - \text{dist}(u) - w(u \rightarrow v)\}$$

(This is defined even when  $\text{dist}(v) = \infty$  or  $\text{dist}(u) = \infty$ , as long as we treat  $\infty$  just like some indescribably large but finite number.) If an edge has zero tension, it's not tense. If we relax an edge  $u \rightarrow v$ , then  $\text{dist}(v)$  decreases  $\text{tension}(u \rightarrow v)$  and  $\text{tension}(u \rightarrow v)$  becomes zero.

Intuitively, we can keep the edges of the graph in some sort of heap, where the key of each edge is its tension. Then we repeatedly pull out the tensest edge  $u \rightarrow v$  and relax it. Then we need to recompute the tension of other edges adjacent to  $v$ . Edges leaving  $v$  possibly become more tense, and edges coming into  $v$  possibly become less tense. So we need a data structure that efficiently supports the operations `INSERT`, `EXTRACTMAX`, `INCREASEKEY`, and `DECREASEKEY`.

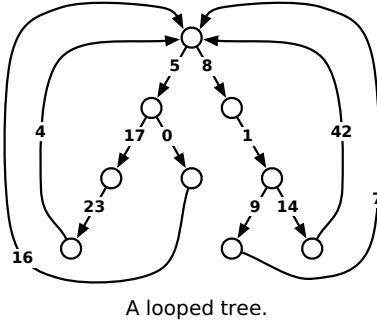
If there are no negative cycles, this algorithm eventually halts with a shortest path tree, but how quickly? Can the same edge be relaxed more than once, and if so, how many times? Is it faster if all the edge weights are positive? Hmm.... This sounds like a good extra credit problem!<sup>8</sup>

## Exercises

1. Prove the following statement for every integer  $i$  and every vertex  $v$ : At the end of the  $i$ th phase of Shimbels algorithm,  $\text{dist}(v)$  is less than or equal to the length of the shortest path  $s \rightsquigarrow v$  consisting of  $i$  or fewer edges.
- \*2. Prove that Ford's generic shortest-path algorithm (while the graph contains a tense edge, relax it) can take exponential time in the worst case when implemented with a stack instead of a priority queue (like Dijkstra) or a queue (like Shimbels algorithm). Specifically, for every positive integer  $n$ , construct a weighted directed  $n$ -vertex graph  $G_n$ , such that the stack-based shortest-path algorithm call `RELAX`  $\Omega(2^n)$  times when  $G_n$  is the input graph. [Hint: Towers of Hanoi.]
- \*3. Prove that Dijkstra's shortest-path algorithm can require exponential time in the worst case when edges are allowed to have negative weight. Specifically, for every positive integer  $n$ , construct a weighted directed  $n$ -vertex graph  $G_n$ , such that Dijkstra's algorithm calls `RELAX`  $\Omega(2^n)$  times when  $G_n$  is the input graph. [Hint: This should be easy if you've already solved the previous problem.]
4. A *looped tree* is a weighted, directed graph built from a binary tree by adding an edge from every leaf back to the root. Every edge has a non-negative weight.

---

<sup>8</sup>I first proposed this bizarre algorithm in 1998, the very first time I taught an algorithms class. As far as I know, nobody has even seriously attempted an analysis. Or maybe it *has* been analyzed, but it requires an exponential number of relaxation steps



- (a) How much time would Dijkstra's algorithm require to compute the shortest path between two vertices  $u$  and  $v$  in a looped tree with  $n$  nodes?
  - (b) Describe and analyze a faster algorithm.
5. (a) Describe a modification of Ford's generic shortest-path algorithm that actually returns a negative cycle if any such cycle is reachable from  $s$ , or a shortest-path tree if there is no such cycle. You may assume that the unmodified algorithm halts in  $O(2^V)$  steps if there is no negative cycle.
- (b) Describe a modification of Shimbel's shortest-path algorithm that actually returns a negative cycle if any such cycle is reachable from  $s$ , or a shortest-path tree if there is no such cycle. The modified algorithm should still run in  $O(VE)$  time.
6. For any edge  $e$  in any graph  $G$ , let  $G \setminus e$  denote the graph obtained by deleting  $e$  from  $G$ .
- (a) Suppose we are given a directed graph  $G$  in which the shortest path  $\sigma$  from vertex  $s$  to vertex  $t$  passes through *every* vertex of  $G$ . Describe an algorithm to compute the shortest-path distance from  $s$  to  $t$  in  $G \setminus e$ , for *every* edge  $e$  of  $G$ , in  $O(E \log V)$  time. Your algorithm should output a set of  $E$  shortest-path distances, one for each edge of the input graph. You may assume that all edge weights are non-negative. [*Hint: If we delete an edge of the original shortest path, how do the old and new shortest paths overlap?*]
  - \*(b) Let  $s$  and  $t$  be *arbitrary* vertices in an *arbitrary* directed graph  $G$ . Describe an algorithm to compute the shortest-path distance from  $s$  to  $t$  in  $G \setminus e$ , for *every* edge  $e$  of  $G$ , in  $O(E \log V)$  time. Again, you may assume that all edge weights are non-negative.
7. Let  $G = (V, E)$  be a connected directed graph with non-negative edge weights, let  $s$  and  $t$  be vertices of  $G$ , and let  $H$  be a subgraph of  $G$  obtained by deleting some edges. Suppose we want to reinsert exactly one edge from  $G$  back into  $H$ , so that the shortest path from  $s$  to  $t$  in the resulting graph is as short as possible. Describe and analyze an algorithm that chooses the best edge to reinsert, in  $O(E \log V)$  time.
8. Negative edges cause problems in shortest-path algorithms because of the possibility of negative cycles. But what if the input graph has no cycles?

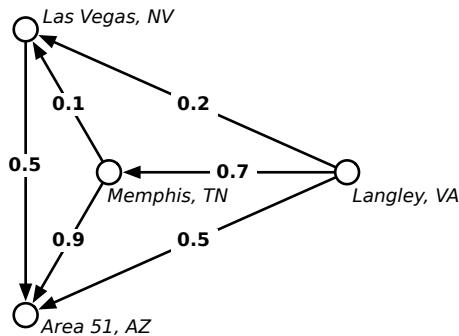
in the worst case, so nobody's ever bothered to publish it.

- (a) Describe an efficient algorithm to compute the shortest path between two nodes  $s$  and  $t$  in a given *directed acyclic graph* with weighted edges. The edge weights could be positive, negative, or zero.
- (b) Describe an efficient algorithm to compute the *longest* path between two nodes  $s$  and  $t$  in a given directed acyclic graph with weighted edges.
9. After a grueling algorithms midterm, you decide to take the bus home. Since you planned ahead, you have a schedule that lists the times and locations of every stop of every bus in Champaign-Urbana. Unfortunately, there isn't a single bus that visits both your exam building and your home; you must transfer between bus lines at least once.
- Describe and analyze an algorithm to determine the sequence of bus rides that will get you home as early as possible, assuming there are  $b$  different bus lines, and each bus stops  $n$  times per day. Your goal is to minimize your *arrival time*, not the time you actually spend traveling. Assume that the buses run exactly on schedule, that you have an accurate watch, and that you are too tired to walk between bus stops.
10. After graduating you accept a job with Aerophobes-Я-Us, the leading traveling agency for people who hate to fly. Your job is to build a system to help customers plan airplane trips from one city to another. All of your customers are afraid of flying (and by extension, airports), so any trip you plan needs to be as short as possible. You know all the departure and arrival times of all the flights on the planet.

Suppose one of your customers wants to fly from city  $X$  to city  $Y$ . Describe an algorithm to find a sequence of flights that minimizes the *total time in transit*—the length of time from the initial departure to the final arrival, including time at intermediate airports waiting for connecting flights. [Hint: Modify the input data and apply Dijkstra's algorithm.]

11. Mulder and Scully have computed, for every road in the United States, the exact probability that someone driving on that road *won't* be abducted by aliens. Agent Mulder needs to drive from Langley, Virginia to Area 51, Nevada. What route should he take so that he has the least chance of being abducted?

More formally, you are given a directed graph  $G = (V, E)$ , where every edge  $e$  has an independent safety probability  $p(e)$ . The *safety* of a path is the product of the safety probabilities of its edges. Design and analyze an algorithm to determine the safest path from a given start vertex  $s$  to a given target vertex  $t$ .



For example, with the probabilities shown above, if Mulder tries to drive directly from Langley to Area 51, he has a 50% chance of getting there without being abducted. If he stops in Memphis, he has a  $0.7 \times 0.9 = 63\%$  chance of arriving safely. If he stops first in Memphis and then in Las Vegas, he has a  $1 - 0.7 \times 0.1 \times 0.5 = 96.5\%$  chance of being abducted! (That's how they got Elvis, you know.)

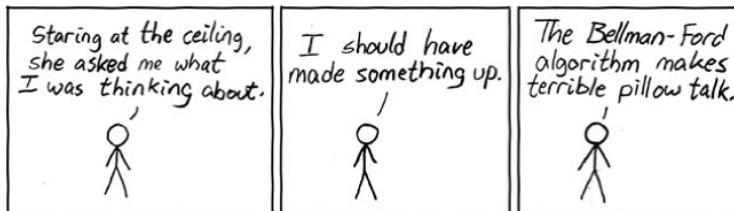
12. On an overnight camping trip in Sunnydale National Park, you are woken from a restless sleep by a scream. As you crawl out of your tent to investigate, a terrified park ranger runs out of the woods, covered in blood and clutching a crumpled piece of paper to his chest. As he reaches your tent, he gasps, “Get out... while... you...”, thrusts the paper into your hands, and falls to the ground. Checking his pulse, you discover that the ranger is stone dead.

You look down at the paper and recognize a map of the park, drawn as an undirected graph, where vertices represent landmarks in the park, and edges represent trails between those landmarks. (Trails start and end at landmarks and do not cross.) You recognize one of the vertices as your current location; several vertices on the boundary of the map are labeled EXIT.

On closer examination, you notice that someone (perhaps the poor dead park ranger) has written a real number between 0 and 1 next to each vertex and each edge. A scrawled note on the back of the map indicates that a number next to an edge is the probability of encountering a vampire along the corresponding trail, and a number next to a vertex is the probability of encountering a vampire at the corresponding landmark. (Vampires can't stand each other's company, so you'll never see more than one vampire on the same trail or at the same landmark.) The note warns you that stepping off the marked trails will result in a slow and painful death.

You glance down at the corpse at your feet. Yes, his death certainly looked painful. Wait, was that a twitch? Are his teeth getting longer? After driving a tent stake through the undead ranger's heart, you wisely decide to leave the park immediately.

Describe and analyze an efficient algorithm to find a path from your current location to an arbitrary EXIT node, such that the total *expected number* of vampires encountered along the path is as small as possible. *Be sure to account for both the vertex probabilities and the edge probabilities!*



— Randall Munroe, xkcd (<http://xkcd.com/69/>)  
Reproduced under a Creative Commons Attribution-NonCommercial 2.5 License

*The tree which fills the arms grew from the tiniest sprout;  
the tower of nine storeys rose from a (small) heap of earth;  
the journey of a thousand li commenced with a single step.*

— Lao-Tzu, *Tao Te Ching*, chapter 64 (6th century BC),  
translated by J. Legge (1891)

*And I would walk five hundred miles,  
And I would walk five hundred more,  
Just to be the man who walks a thousand miles  
To fall down at your door.*

— The Proclaimers, “Five Hundred Miles (I’m Gonna Be)”,  
*Sunshine on Leith* (2001)

*Almost there... Almost there...*

— Red Leader [Drewe Henley], *Star Wars* (1977)

## 20 All-Pairs Shortest Paths

### 20.1 The Problem

In the previous lecture, we saw algorithms to find the shortest path from a source vertex  $s$  to a target vertex  $t$  in a directed graph. As it turns out, the best algorithms for this problem actually find the shortest path from  $s$  to every possible target (or from every possible source to  $t$ ) by constructing a shortest path tree. The shortest path tree specifies two pieces of information for each node  $v$  in the graph:

- $dist(v)$  is the length of the shortest path (if any) from  $s$  to  $v$ ;
- $pred(v)$  is the second-to-last vertex (if any) the shortest path (if any) from  $s$  to  $v$ .

In this lecture, we want to generalize the shortest path problem even further. In the *all pairs shortest path* problem, we want to find the shortest path from *every* possible source to *every* possible destination. Specifically, for every pair of vertices  $u$  and  $v$ , we need to compute the following information:

- $dist(u, v)$  is the length of the shortest path (if any) from  $u$  to  $v$ ;
- $pred(u, v)$  is the second-to-last vertex (if any) on the shortest path (if any) from  $u$  to  $v$ .

For example, for any vertex  $v$ , we have  $dist(v, v) = 0$  and  $pred(v, v) = \text{NULL}$ . If the shortest path from  $u$  to  $v$  is only one edge long, then  $dist(u, v) = w(u \rightarrow v)$  and  $pred(u, v) = u$ . If there is *no* shortest path from  $u$  to  $v$ —either because there’s no path at all, or because there’s a negative cycle—then  $dist(u, v) = \infty$  and  $pred(v, v) = \text{NULL}$ .

The output of our shortest path algorithms will be a pair of  $V \times V$  arrays encoding all  $V^2$  distances and predecessors. Many maps include a distance matrix—to find the distance from (say) Champaign to (say) Columbus, you would look in the row labeled ‘Champaign’ and the column labeled ‘Columbus’. In these notes, I’ll focus almost exclusively on computing the distance array. The predecessor array, from which you would compute the actual shortest paths, can be computed with only minor additions to the algorithms I’ll describe (hint, hint).

### 20.2 Lots of Single Sources

The obvious solution to the all-pairs shortest path problem is just to run a single-source shortest path algorithm  $V$  times, once for every possible source vertex! Specifically, to fill in the one-dimensional subarray  $dist[s, \cdot]$ , we invoke either Dijkstra’s or Shimbels algorithm starting at the source vertex  $s$ .

**OBVIOUSAPSP( $V, E, w$ ):**

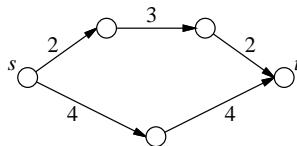
for every vertex  $s$

$dist[s, \cdot] \leftarrow \text{SSSP}(V, E, w, s)$

The running time of this algorithm depends on which single-source shortest path algorithm we use. If we use Shimbels algorithm, the overall running time is  $\Theta(V^2E) = O(V^4)$ . If all the edge weights are non-negative, we can use Dijkstra's algorithm instead, which decreases the running time to  $\Theta(VE + V^2 \log V) = O(V^3)$ . For graphs with negative edge weights, Dijkstra's algorithm can take exponential time, so we can't get this improvement directly.

### 20.3 Reweighting

One idea that occurs to most people is increasing the weights of all the edges by the same amount so that all the weights become positive, and then applying Dijkstra's algorithm. Unfortunately, this simple idea doesn't work. Different paths change by different amounts, which means the shortest paths in the reweighted graph may not be the same as in the original graph.



Increasing all the edge weights by 2 changes the shortest path  $s$  to  $t$ .

However, there is a more complicated method for reweighting the edges in a graph. Suppose each vertex  $v$  has some associated cost  $c(v)$ , which might be positive, negative, or zero. We can define a new weight function  $w'$  as follows:

$$w'(u \rightarrow v) = c(u) + w(u \rightarrow v) - c(v)$$

To give some intuition, imagine that when we leave vertex  $u$ , we have to pay an exit tax of  $c(u)$ , and when we enter  $v$ , we get  $c(v)$  as an entrance gift.

Now it's not too hard to show that the shortest paths with the new weight function  $w'$  are exactly the same as the shortest paths with the original weight function  $w$ . In fact, for any path  $u \rightsquigarrow v$  from one vertex  $u$  to another vertex  $v$ , we have

$$w'(u \rightsquigarrow v) = c(u) + w(u \rightsquigarrow v) - c(v).$$

We pay  $c(u)$  in exit fees, plus the original weight of the path, minus the  $c(v)$  entrance gift. At every intermediate vertex  $x$  on the path, we get  $c(x)$  as an entrance gift, but then immediately pay it back as an exit tax!

### 20.4 Johnson's Algorithm

Johnson's all-pairs shortest path algorithm finds a cost  $c(v)$  for each vertex, so that when the graph is reweighted, every edge has non-negative weight.

Suppose the graph has a vertex  $s$  that has a path to every other vertex. Johnson's algorithm computes the shortest paths from  $s$  to every other vertex, using Shimbels algorithm (which doesn't care if the edge weights are negative), and then sets  $c(v) \leftarrow dist(s, v)$ , so the new weight of every edge is

$$w'(u \rightarrow v) = dist(s, u) + w(u \rightarrow v) - dist(s, v).$$

Why are all these new weights non-negative? Because otherwise, Shimbels algorithm wouldn't be finished! Recall that an edge  $u \rightarrow v$  is *tense* if  $dist(s, u) + w(u \rightarrow v) < dist(s, v)$ , and that single-source shortest path algorithms eliminate all tense edges. The only exception is if the graph has a negative cycle, but then shortest paths aren't defined, and Johnson's algorithm simply aborts.

But what if the graph *doesn't* have a vertex  $s$  that can reach everything? No matter where we start Shimbels algorithm, some of those vertex costs will be infinite. Johnson's algorithm avoids this problem by adding a new vertex  $s$  to the graph, with zero-weight edges going from  $s$  to every other vertex, but no edges going back into  $s$ . This addition doesn't change the shortest paths between any other pair of vertices, because there are no paths into  $s$ .

So here's Johnson's algorithm in all its glory.

```
JOHNSONAPSP( $V, E, w$ ) :
  create a new vertex  $s$ 
  for every vertex  $v$ 
     $w(s \rightarrow v) \leftarrow 0$ 
     $w(v \rightarrow s) \leftarrow \infty$ 
   $dist[s, \cdot] \leftarrow \text{SHIMBEL}(V, E, w, s)$ 
  abort if SHIMBEL found a negative cycle
  for every edge  $(u, v) \in E$ 
     $w'(u \rightarrow v) \leftarrow dist[s, u] + w(u \rightarrow v) - dist[s, v]$ 
  for every vertex  $u$ 
     $dist[u, \cdot] \leftarrow \text{DIJKSTRA}(V, E, w', u)$ 
    for every vertex  $v$ 
       $dist[u, v] \leftarrow dist[u, v] - dist[s, u] + dist[s, v]$ 
```

The algorithm spends  $\Theta(V)$  time adding the artificial start vertex  $s$ ,  $\Theta(VE)$  time running SHIMBEL,  $O(E)$  time reweighting the graph, and then  $\Theta(VE + V^2 \log V)$  running  $V$  passes of Dijkstra's algorithm. Thus, the overall running time is  $\Theta(VE + V^2 \log V)$ .

## 20.5 Dynamic Programming

There's a completely different solution to the all-pairs shortest path problem that uses dynamic programming instead of a single-source algorithm. For *dense* graphs where  $E = \Omega(V^2)$ , the dynamic programming approach eventually leads to the same  $O(V^3)$  running time as Johnson's algorithm, but with a much simpler algorithm. In particular, the new algorithm avoids Dijkstra's algorithm, which gets its efficiency from Fibonacci heaps, which are rather easy to screw up in the implementation. **In the rest of this lecture, I will assume that the input graph contains no negative cycles.**

As usual for dynamic programming algorithms, we first need to come up with a recursive formulation of the problem. Here is an "obvious" recursive definition for  $dist(u, v)$ :

$$dist(u, v) = \begin{cases} 0 & \text{if } u = v \\ \min_x (dist(u, x) + w(x \rightarrow v)) & \text{otherwise} \end{cases}$$

In other words, to find the shortest path from  $u$  to  $v$ , try all possible predecessors  $x$ , compute the shortest path from  $u$  to  $x$ , and then add the last edge  $u \rightarrow v$ . **Unfortunately, this recurrence doesn't work!** In order to compute  $dist(u, v)$ , we first have to compute  $dist(u, x)$  for every other vertex  $x$ , but to compute any  $dist(u, x)$ , we first need to compute  $dist(u, v)$ . We're stuck in an infinite loop!

To avoid this circular dependency, we need an additional parameter that decreases at each recursion, eventually reaching zero at the base case. One possibility is to include the number of edges in the

shortest path as this third magic parameter. So let's define  $dist(u, v, k)$  to be the length of the shortest path from  $u$  to  $v$  that uses *at most*  $k$  edges. Since we know that the shortest path between any two vertices has at most  $V - 1$  vertices, what we're really trying to compute is  $dist(u, v, V - 1)$ .

After a little thought, we get the following recurrence.

$$dist(u, v, k) = \begin{cases} 0 & \text{if } u = v \\ \infty & \text{if } k = 0 \text{ and } u \neq v \\ \min_x (dist(u, x, k - 1) + w(x \rightarrow v)) & \text{otherwise} \end{cases}$$

Just like last time, the recurrence tries all possible predecessors of  $v$  in the shortest path, but now the recursion actually bottoms out when  $k = 0$ .

Now it's not difficult to turn this recurrence into a dynamic programming algorithm. Even before we write down the algorithm, though, we can tell that its running time will be  $\Theta(V^4)$  simply because recurrence has four variables— $u$ ,  $v$ ,  $k$ , and  $x$ —each of which can take on  $V$  different values. Except for the base cases, the algorithm itself is just four nested for loops. To make the algorithm a little shorter, let's assume that  $w(v \rightarrow v) = 0$  for every vertex  $v$ .

```
DYNAMICPROGRAMMINGAPSP( $V, E, w$ ):  
    for all vertices  $u$   
        for all vertices  $v$   
            if  $u = v$   
                 $dist[u, v, 0] \leftarrow 0$   
            else  
                 $dist[u, v, 0] \leftarrow \infty$   
    for  $k \leftarrow 1$  to  $V - 1$   
        for all vertices  $u$   
            for all vertices  $v$   
                 $dist[u, v, k] \leftarrow \infty$   
                for all vertices  $x$   
                    if  $dist[u, v, k] > dist[u, x, k - 1] + w(x \rightarrow v)$   
                         $dist[u, v, k] \leftarrow dist[u, x, k - 1] + w(x \rightarrow v)$ 
```

In fact, this algorithm is almost the same as running Shimbel's algorithm once from each source vertex. The main difference is the innermost loop, which in Shimbel's algorithm would read “for all edges  $x \rightarrow v$ ”. This simple change improves the running time to  $\Theta(V^2E)$ , assuming the graph is stored in an adjacency list. Other differences are just as in the dynamic programming development of Shimbel's algorithm—we don't need the inner loop over vertices  $v$ , and we only need a two-dimensional table. After the  $k$ th iteration of the main loop in the following algorithm, the value of  $dist[u, v]$  lies between the true shortest path distance from  $u$  to  $v$  and the value  $dist[u, v, k]$  computed in the previous algorithm.

```
SHIMBELAPSP( $V, E, w$ ):  
    for all vertices  $u$   
        for all vertices  $v$   
            if  $u = v$   
                 $dist[u, v] \leftarrow 0$   
            else  
                 $dist[u, v] \leftarrow \infty$   
    for  $k \leftarrow 1$  to  $V - 1$   
        for all vertices  $u$   
            for all edges  $x \rightarrow v$   
                if  $dist[u, v] > dist[u, x] + w(x \rightarrow v)$   
                     $dist[u, v] \leftarrow dist[u, x] + w(x \rightarrow v)$ 
```

## 20.6 Divide and Conquer

But we can make a more significant improvement. The recurrence we just used broke the shortest path into a slightly shorter path and a single edge, by considering all predecessors. Instead, let's break it into two shorter paths at the middle vertex on the path. This idea gives us a different recurrence for  $\text{dist}(u, v, k)$ . Once again, to simplify things, let's assume  $w(v \rightarrow v) = 0$ .

$$\text{dist}(u, v, k) = \begin{cases} w(u \rightarrow v) & \text{if } k = 1 \\ \min_x (\text{dist}(u, x, k/2) + \text{dist}(x, v, k/2)) & \text{otherwise} \end{cases}$$

This recurrence only works when  $k$  is a power of two, since otherwise we might try to find the shortest path with a fractional number of edges! But that's not really a problem, since  $\text{dist}(u, v, 2^{\lceil \lg V \rceil})$  gives us the overall shortest distance from  $u$  to  $v$ . Notice that we use the base case  $k = 1$  instead of  $k = 0$ , since we can't use half an edge.

Once again, a dynamic programming solution is straightforward. Even before we write down the algorithm, we can tell the running time is  $\Theta(V^3 \log V)$ —we consider  $V$  possible values of  $u$ ,  $v$ , and  $x$ , but only  $\lceil \lg V \rceil$  possible values of  $k$ .

```
FASTDYNAMICPROGRAMMINGAPSP( $V, E, w$ ):
    for all vertices  $u$ 
        for all vertices  $v$ 
             $\text{dist}[u, v, 0] \leftarrow w(u \rightarrow v)$ 
    for  $i \leftarrow 1$  to  $\lceil \lg V \rceil$        $\langle\langle k = 2^i \rangle\rangle$ 
        for all vertices  $u$ 
            for all vertices  $v$ 
                 $\text{dist}[u, v, i] \leftarrow \infty$ 
                for all vertices  $x$ 
                    if  $\text{dist}[u, v, i] > \text{dist}[u, x, i - 1] + \text{dist}[x, v, i - 1]$ 
                         $\text{dist}[u, v, i] \leftarrow \text{dist}[u, x, i - 1] + \text{dist}[x, v, i - 1]$ 
```

This algorithm is *not* the same as  $V$  invocations of any single-source algorithm; we can't replace the innermost loop over vertices with a loop over edges. However, we can remove the last dimension of the table, using  $\text{dist}[u, v]$  everywhere in place of  $\text{dist}[u, v, i]$ , just as in Shimbels algorithm, thereby reducing the space from  $O(V^3)$  to  $O(V^2)$ .

```
FASTSHIMBELAPSP( $V, E, w$ ):
    for all vertices  $u$ 
        for all vertices  $v$ 
             $\text{dist}[u, v] \leftarrow w(u \rightarrow v)$ 
    for  $i \leftarrow 1$  to  $\lceil \lg V \rceil$ 
        for all vertices  $u$ 
            for all vertices  $v$ 
                for all vertices  $x$ 
                    if  $\text{dist}[u, v] > \text{dist}[u, x] + \text{dist}[x, v]$ 
                         $\text{dist}[u, v] \leftarrow \text{dist}[u, x] + \text{dist}[x, v]$ 
```

## 20.7 Aside: ‘Funny’ Matrix Multiplication

There is a very close connection (first observed by Shimbels, and later independently by Bellman) between computing shortest paths in a directed graph and computing powers of a square matrix. Compare the following algorithm for multiplying two  $n \times n$  matrices  $A$  and  $B$  with the inner loop of our first dynamic programming algorithm. (I've changed the variable names in the second algorithm slightly to make the similarity clearer.)

```
MATRIXMULTIPLY( $A, B$ ):
  for  $i \leftarrow 1$  to  $n$ 
    for  $j \leftarrow 1$  to  $n$ 
       $C[i, j] \leftarrow 0$ 
      for  $k \leftarrow 1$  to  $n$ 
         $C[i, j] \leftarrow C[i, j] + A[i, k] \cdot B[k, j]$ 
```

```
APSPINNERLOOP:
  for all vertices  $u$ 
    for all vertices  $v$ 
       $D'[u, v] \leftarrow \infty$ 
      for all vertices  $x$ 
         $D'[u, v] \leftarrow \min \{D'[u, v], D[u, x] + w[x, v]\}$ 
```

The *only* difference between these two algorithms is that we use addition instead of multiplication and minimization instead of addition. For this reason, the shortest path inner loop is often referred to as ‘funny’ matrix multiplication.

**DYNAMICPROGRAMMINGAPSP** is the standard iterative algorithm for computing the  $(V - 1)$ th ‘funny power’ of the weight matrix  $w$ . The first set of for loops sets up the ‘funny identity matrix’, with zeros on the main diagonal and infinity everywhere else. Then each iteration of the second main for loop computes the next ‘funny power’. **FASTDYNAMICPROGRAMMINGAPSP** replaces this iterative method for computing powers with repeated squaring, exactly like we saw at the beginning of the semester. The fast algorithm is simplified slightly by the fact that unless there are negative cycles, every ‘funny power’ after the  $V$ th is the same.

There are faster methods for multiplying matrices, similar to Karatsuba’s divide-and-conquer algorithm for multiplying integers. (Google for ‘Strassen’s algorithm’.) Unfortunately, these algorithms use subtraction, and there’s no ‘funny’ equivalent of subtraction. (What’s the inverse operation for min?) So at least for general graphs, there seems to be no way to speed up the inner loop of our dynamic programming algorithms.

Fortunately, this isn’t true. There is a beautiful randomized algorithm, due to Noga Alon, Zvi Galil, Oded Margalit\*, and Moni Naor,<sup>1</sup> that computes all-pairs shortest paths in undirected graphs in  $O(M(V) \log^2 V)$  expected time, where  $M(V)$  is the time to multiply two  $V \times V$  integer matrices. A simplified version of this algorithm for *unweighted* graphs was discovered by Raimund Seidel.<sup>2</sup>

## 20.8 Floyd and Warshall’s Algorithm

Our fast dynamic programming algorithm is still a factor of  $O(\log V)$  slower than Johnson’s algorithm. A different formulation due to Floyd and Warshall removes this logarithmic factor. Their insight was to use a different third parameter in the recurrence.

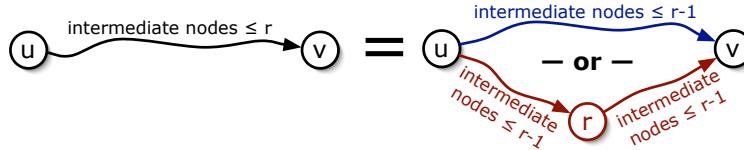
Number the vertices arbitrarily from 1 to  $V$ . For every pair of vertices  $u$  and  $v$  and every integer  $r$ , we define a path  $\pi(u, v, r)$  as follows:

$\pi(u, v, r) :=$  the shortest path from  $u$  to  $v$  where every intermediate vertex (that is, every vertex except  $u$  and  $v$ ) is numbered at most  $r$ .

<sup>1</sup>N. Alon, Z. Galil, O. Margalit\*, and M. Naor. Witnesses for Boolean matrix multiplication and for shortest paths. *Proc. 33rd FOCS* 417-426, 1992. See also N. Alon, Z. Galil, O. Margalit\*. On the exponent of the all pairs shortest path problem. *Journal of Computer and System Sciences* 54(2):255–262, 1997.

<sup>2</sup>R. Seidel. On the all-pairs-shortest-path problem in unweighted undirected graphs. *Journal of Computer and System Sciences*, 51(3):400-403, 1995. This is one of the few algorithms papers where (in the conference version at least) the algorithm is completely described and analyzed *in the abstract* of the paper.

If  $r = 0$ , we aren't allowed to use any intermediate vertices, so  $\pi(u, v, 0)$  is just the edge (if any) from  $u$  to  $v$ . If  $r > 0$ , then either  $\pi(u, v, r)$  goes through the vertex numbered  $r$ , or it doesn't. If  $\pi(u, v, r)$  does contain vertex  $r$ , it splits into a subpath from  $u$  to  $r$  and a subpath from  $r$  to  $v$ , where every intermediate vertex in these two subpaths is numbered at most  $r - 1$ . Moreover, the subpaths are as short as possible with this restriction, so they must be  $\pi(u, r, r - 1)$  and  $\pi(r, v, r - 1)$ . On the other hand, if  $\pi(u, v, r)$  does not go through vertex  $r$ , then every intermediate vertex in  $\pi(u, v, r)$  is numbered at most  $r - 1$ ; since  $\pi(u, v, r)$  must be the *shortest* such path, we have  $\pi(u, v, r) = \pi(u, v, r - 1)$ .



Recursive structure of the restricted shortest path  $\pi(u, v, r)$ .

This recursive structure implies the following recurrence for the length of  $\pi(u, v, r)$ , which we will denote by  $dist(u, v, r)$ :

$$dist(u, v, r) = \begin{cases} w(u \rightarrow v) & \text{if } r = 0 \\ \min \{ dist(u, v, r - 1), dist(u, r, r - 1) + dist(r, v, r - 1) \} & \text{otherwise} \end{cases}$$

We need to compute the shortest path distance from  $u$  to  $v$  with no restrictions, which is just  $dist(u, v, V)$ . Once again, we should immediately see that a dynamic programming algorithm will implement this recurrence in  $\Theta(V^3)$  time.

```
FLOYDWARSHALL( $V, E, w$ ):
  for all vertices  $u$ 
    for all vertices  $v$ 
       $dist[u, v, 0] \leftarrow w(u \rightarrow v)$ 
  for  $r \leftarrow 1$  to  $V$ 
    for all vertices  $u$ 
      for all vertices  $v$ 
        if  $dist[u, v, r - 1] < dist[u, r, r - 1] + dist[r, v, r - 1]$ 
           $dist[u, v, r] \leftarrow dist[u, v, r - 1]$ 
        else
           $dist[u, v, r] \leftarrow dist[u, r, r - 1] + dist[r, v, r - 1]$ 
```

Just like our earlier algorithms, we can simplify the algorithm by removing the third dimension of the memoization table. Also, because the vertex numbering was chosen arbitrary, so there's no reason to refer to it explicitly in the pseudocode.

```
FLOYDWARSHALL2( $V, E, w$ ):
  for all vertices  $u$ 
    for all vertices  $v$ 
       $dist[u, v] \leftarrow w(u \rightarrow v)$ 
  for all vertices  $r$ 
    for all vertices  $u$ 
      for all vertices  $v$ 
        if  $dist[u, v] > dist[u, r] + dist[r, v]$ 
           $dist[u, v] \leftarrow dist[u, r] + dist[r, v]$ 
```

Now compare this algorithm with FASTSHIMBELAPSP. Instead of  $O(\log V)$  passes through all triples of vertices, FLOYDWARSHALL2 only requires a single pass, but only because it uses a different nesting order for the three for-loops!

## Exercises

1. Let  $G = (V, E)$  be a directed graph with weighted edges; edge weights could be positive, negative, or zero.
  - (a) How could we delete some node  $v$  from this graph, without changing the shortest-path distance between any other pair of nodes? Describe an algorithm that constructs a directed graph  $G' = (V', E')$  with weighted edges, where  $V' = V \setminus \{v\}$ , and the shortest-path distance between any two nodes in  $H$  is equal to the shortest-path distance between the same two nodes in  $G$ , in  $O(V^2)$  time.
  - (b) Now suppose we have already computed all shortest-path distances in  $G'$ . Describe an algorithm to compute the shortest-path distances from  $v$  to every other node, and from every other node to  $v$ , in the original graph  $G$ , in  $O(V^2)$  time.
  - (c) Combine parts (a) and (b) into another all-pairs shortest path algorithm that runs in  $O(V^3)$  time. (The resulting algorithm is *not* the same as Floyd-Warshall!)
2. All of the algorithms discussed in this lecture fail if the graph contains a negative cycle. Johnson's algorithm detects the negative cycle in the initialization phase (via Shimbels algorithm) and aborts; the dynamic programming algorithms just return incorrect results. However, all of these algorithms can be modified to return correct shortest-path distances, even in the presence of negative cycles. Specifically, if there is a path from vertex  $u$  to a negative cycle and a path from that negative cycle to vertex  $v$ , the algorithm should report that  $\text{dist}[u, v] = -\infty$ . If there is no directed path from  $u$  to  $v$ , the algorithm should return  $\text{dist}[u, v] = \infty$ . Otherwise,  $\text{dist}[u, v]$  should equal the length of the shortest directed path from  $u$  to  $v$ .
  - (a) Describe how to modify Johnson's algorithm to return the correct shortest-path distances, even if the graph has negative cycles.
  - (b) Describe how to modify the Floyd-Warshall algorithm (FLOYDWARSHALL2) to return the correct shortest-path distances, even if the graph has negative cycles.
3. All of the shortest-path algorithms described in this note can also be modified to return an explicit description of some negative cycle, instead of simply reporting that a negative cycle exists.
  - (a) Describe how to modify Johnson's algorithm to return either the matrix of shortest-path distances or a negative cycle.
  - (b) Describe how to modify the Floyd-Warshall algorithm (FLOYDWARSHALL2) to return either the matrix of shortest-path distances or a negative cycle.

If the graph contains more than one negative cycle, your algorithms may choose one arbitrarily.

4. Let  $G = (V, E)$  be a directed graph with weighted edges; edge weights could be positive, negative, or zero. Suppose the vertices of  $G$  are partitioned into  $k$  disjoint subsets  $V_1, V_2, \dots, V_k$ ; that is, every vertex of  $G$  belongs to exactly one subset  $V_i$ . For each  $i$  and  $j$ , let  $\delta(i, j)$  denote the minimum shortest-path distance between vertices in  $V_i$  and vertices in  $V_j$ :

$$\delta(i, j) = \min \{ \text{dist}(u, v) \mid u \in V_i \text{ and } v \in V_j \}.$$

Describe an algorithm to compute  $\delta(i, j)$  for all  $i$  and  $j$  in time  $O(V^2 + kE \log E)$ .

5. Recall<sup>3</sup> that a deterministic finite automaton (DFA) is formally defined as a tuple  $M = (\Sigma, Q, q_0, F, \delta)$ , where the finite set  $\Sigma$  is the input alphabet, the finite set  $Q$  is the set of states,  $q_0 \in Q$  is the start state,  $F \subseteq Q$  is the set of final (accepting) states, and  $\delta: Q \times \Sigma \rightarrow Q$  is the transition function. Equivalently, a DFA is a directed (multi-)graph with labeled edges, such that each symbol in  $\Sigma$  is the label of exactly one edge leaving any vertex. There is a special ‘start’ vertex  $q_0$ , and a subset of the vertices are marked as ‘accepting states’. Any string in  $\Sigma^*$  describes a unique walk starting at  $q_0$ ; a string in  $\Sigma^*$  is *accepted* by  $M$  if this walk ends at a vertex in  $F$ .

Stephen Kleene<sup>4</sup> proved that the language accepted by any DFA is identical to the language described by some regular expression. This problem asks you to develop a variant of the Floyd-Warshall all-pairs shortest path algorithm that computes a regular expression that is equivalent to the language accepted by a given DFA.

Suppose the input DFA  $M$  has  $n$  states, numbered from 1 to  $n$ , where (without loss of generality) the start state is state 1. Let  $L(i, j, r)$  denote the set of all words that describe walks in  $M$  from state  $i$  to state  $j$ , where every intermediate state lies in the subset  $\{1, 2, \dots, r\}$ ; thus, the language accepted by the DFA is exactly

$$\bigcup_{q \in F} L(1, q, n).$$

Let  $R(i, j, r)$  be a regular expression that describes the language  $L(i, j, r)$ .

- (a) What is the regular expression  $R(i, j, 0)$ ?
  - (b) Write a recurrence for the regular expression  $R(i, j, r)$  in terms of regular expressions of the form  $R(i', j', r - 1)$ .
  - (c) Describe a polynomial-time algorithm to compute  $R(i, j, n)$  for all states  $i$  and  $j$ . (Assume that you can concatenate two regular expressions in  $O(1)$  time.)
- \*6. Let  $G = (V, E)$  be an undirected, unweighted, connected,  $n$ -vertex graph, represented by the adjacency matrix  $A[1..n, 1..n]$ . In this problem, we will derive Seidel’s sub-cubic algorithm to compute the  $n \times n$  matrix  $D[1..n, 1..n]$  of shortest-path distances using fast matrix multiplication. Assume that we have a subroutine `MATRIXMULTIPLY` that multiplies two  $n \times n$  matrices in  $\Theta(n^\omega)$  time, for some unknown constant  $\omega \geq 2$ .<sup>5</sup>

---

<sup>3</sup>Automata theory is a prerequisite for the undergraduate algorithms class at UIUC.

<sup>4</sup>Pronounced ‘clay knee’, not ‘clean’ or ‘clean-ee’ or ‘clay-nuh’ or ‘dimaggio’.

<sup>5</sup>The matrix multiplication algorithm you already know runs in  $\Theta(n^3)$  time, but this is not the fastest algorithm known. The current record is  $\omega \approx 2.376$ , due to Don Coppersmith and Shmuel Winograd. Determining the smallest possible value of  $\omega$  is a long-standing open problem; many people believe there is an undiscovered  $O(n^2)$ -time algorithm for matrix multiplication.

- (a) Let  $G^2$  denote the graph with the same vertices as  $G$ , where two vertices are connected by a edge if and only if they are connected by a path of length at most 2 in  $G$ . Describe an algorithm to compute the adjacency matrix of  $G^2$  using a single call to MATRIXMULTIPLY and  $O(n^2)$  additional time.
- (b) Suppose we discover that  $G^2$  is a complete graph. Describe an algorithm to compute the matrix  $D$  of shortest path distances in  $O(n^2)$  additional time.
- (c) Let  $D^2$  denote the (recursively computed) matrix of shortest-path distances in  $G^2$ . Prove that the shortest-path distance from node  $i$  to node  $j$  is either  $2 \cdot D^2[i, j]$  or  $2 \cdot D^2[i, j] - 1$ .
- (d) Suppose  $G^2$  is not a complete graph. Let  $X = D^2 \cdot A$ , and let  $\deg(i)$  denote the degree of vertex  $i$  in the original graph  $G$ . Prove that the shortest-path distance from node  $i$  to node  $j$  is  $2 \cdot D^2[i, j]$  if and only if  $X[i, j] \geq D^2[i, j] \cdot \deg(i)$ .
- (e) Describe an algorithm to compute the matrix of shortest-path distances in  $G$  in  $O(n^\omega \log n)$  time.

**Col. Hogan:** One of these wires disconnects the fuse,  
the other one fires the bomb.

Which one would you cut, Shultz?

**Sgt. Schultz:** Don't ask me, this is a decision for an officer.

**Col. Hogan:** All right. Which wire, Colonel Klink?

**Col. Klink:** This one. [points to the white wire]

**Col. Hogan:** You're sure?

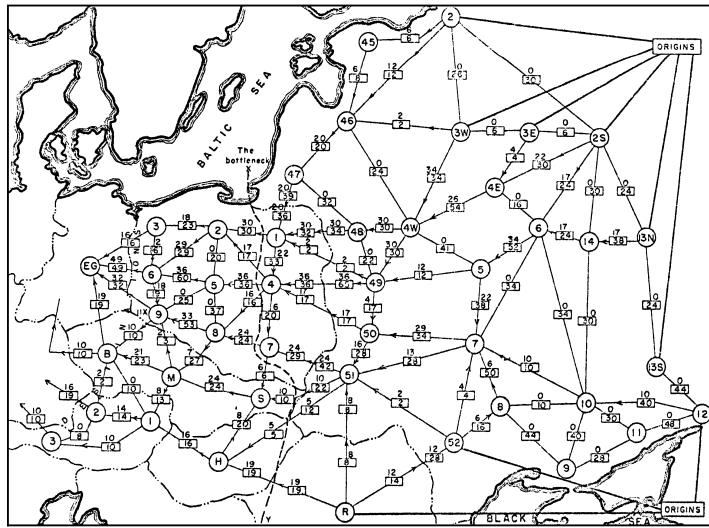
**Col. Klink:** Yes.

[Hogan cuts the black wire; the bomb stops ticking]

**Col. Klink:** If you knew which wire it was, why did you ask me?

**Col. Hogan:** I wasn't sure which was the right one,  
but I was certain you'd pick the wrong one.

— "A Klink, a Bomb, and a Short Fuse", Hogan's Heroes (1966)



Harris and Ross's map of the Warsaw Pact rail network

## 21 Maximum Flows and Minimum Cuts

In the mid-1950s, Air Force researchers T. E. Harris and F. S. Ross published a classified report studying the rail network that linked the Soviet Union to its satellite countries in Eastern Europe. The network was modeled as a graph with 44 vertices, representing geographic regions, and 105 edges, representing links between those regions in the rail network. Each edge was given a weight, representing the rate at which material could be shipped from one region to the next. Essentially by trial and error, they determined both the maximum amount of stuff that could be moved from Russia into Europe, as well as the cheapest way to disrupt the network by removing links (or in less abstract terms, blowing up train tracks), which they called 'the bottleneck'. Their results (including the figure at the top of the page) were only declassified in 1999.<sup>1</sup>

This one of the first recorded applications of the *maximum flow* and *minimum cut* problems. For both problems, the input is a directed graph  $G = (V, E)$ , along with special vertices  $s$  and  $t$  called the *source* and *target*. As in the previous lectures, I will use  $u \rightarrow v$  to denote the directed edge from vertex  $u$

<sup>1</sup>Both the map and the story were taken from Alexander Schrijver's fascinating survey 'On the history of combinatorial optimization (till 1960)'.

to vertex  $v$ . Intuitively, the maximum flow problem asks for the largest amount of material that can be transported from one vertex to another; the minimum cut problem asks for the minimum damage needed to separate two vertices.

## 21.1 Flows

An  $(s, t)$ -flow (or just a *flow* if the source and target are clear from context) is a function  $f : E \rightarrow \mathbb{R}_{\geq 0}$  that satisfies the following **conservation constraint** at every vertex  $v$  except possibly  $s$  and  $t$ :

$$\sum_u f(u \rightarrow v) = \sum_w f(v \rightarrow w).$$

In English, the total flow into  $v$  is equal to the total flow out of  $v$ . To keep the notation simple, we define  $f(u \rightarrow v) = 0$  if there is no edge  $u \rightarrow v$  in the graph. The **value** of the flow  $f$ , denoted  $|f|$ , is the total net flow out of the source vertex  $s$ :

$$|f| := \sum_w f(s \rightarrow w) - \sum_u f(u \rightarrow s).$$

It's not hard to prove that  $|f|$  is also equal to the total net flow *into* the target vertex  $t$ , as follows. To simplify notation, let  $\partial f(v)$  denote the total net flow out of any vertex  $v$ :

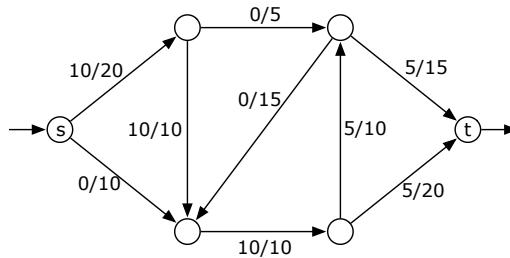
$$\partial f(v) := \sum_u f(u \rightarrow v) - \sum_w f(v \rightarrow w).$$

The conservation constraint implies that  $\partial f(v) = 0$  for every vertex  $v$  except  $s$  and  $t$ , so

$$\sum_v \partial f(v) = \partial f(s) + \partial f(t).$$

On the other hand, any flow that leaves one vertex must enter another vertex, so we must have  $\sum_v \partial f(v) = 0$ . It follows immediately that  $|f| = \partial f(s) = -\partial f(t)$ .

Now suppose we have another function  $c : E \rightarrow \mathbb{R}_{\geq 0}$  that assigns a non-negative **capacity**  $c(e)$  to each edge  $e$ . We say that a flow  $f$  is **feasible** (with respect to  $c$ ) if  $f(e) \leq c(e)$  for every edge  $e$ . Most of the time we will consider only flows that are feasible with respect to some fixed capacity function  $c$ . We say that a flow  $f$  **saturates** edge  $e$  if  $f(e) = c(e)$ , and **avoids** edge  $e$  if  $f(e) = 0$ . The **maximum flow problem** is to compute a feasible  $(s, t)$ -flow in a given directed graph, with a given capacity function, whose value is as large as possible.



An  $(s, t)$ -flow with value 10. Each edge is labeled with its flow/capacity.

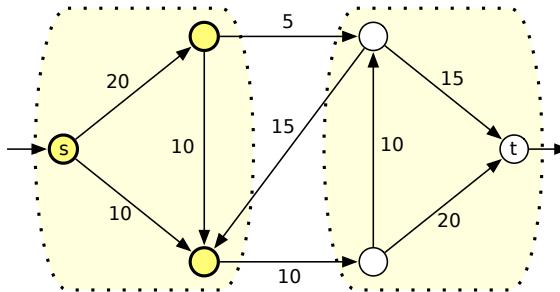
## 21.2 Cuts

An  $(s, t)$ -cut (or just *cut* if the source and target are clear from context) is a partition of the vertices into disjoint subsets  $S$  and  $T$ —meaning  $S \cup T = V$  and  $S \cap T = \emptyset$ —where  $s \in S$  and  $t \in T$ .

If we have a capacity function  $c: E \rightarrow \mathbb{R}_{\geq 0}$ , the *capacity* of a cut is the sum of the capacities of the edges that start in  $S$  and end in  $T$ :

$$\|S, T\| := \sum_{v \in S} \sum_{w \in T} c(v \rightarrow w).$$

(Again, if  $v \rightarrow w$  is not an edge in the graph, we assume  $c(v \rightarrow w) = 0$ .) Notice that the definition is asymmetric; edges that start in  $T$  and end in  $S$  are unimportant. The *minimum cut problem* is to compute an  $(s, t)$ -cut whose capacity is as large as possible.



An  $(s, t)$ -cut with capacity 15. Each edge is labeled with its capacity.

Intuitively, the minimum cut is the cheapest way to disrupt all flow from  $s$  to  $t$ . Indeed, it is not hard to show that **the value of any feasible  $(s, t)$ -flow is at most the capacity of any  $(s, t)$ -cut**. Choose your favorite flow  $f$  and your favorite cut  $(S, T)$ , and then follow the bouncing inequalities:

$$\begin{aligned}
 |f| &= \sum_w f(s \rightarrow w) - \sum_u f(u \rightarrow s) && \text{by definition} \\
 &= \sum_{v \in S} \left( \sum_w f(v \rightarrow w) - \sum_u f(u \rightarrow v) \right) && \text{by the conservation constraint} \\
 &= \sum_{v \in S} \left( \sum_{w \in T} f(v \rightarrow w) - \sum_{u \in T} f(u \rightarrow v) \right) && \text{removing duplicate edges} \\
 &\leq \sum_{v \in S} \sum_{w \in T} f(v \rightarrow w) && \text{since } f(u \rightarrow v) \geq 0 \\
 &\leq \sum_{v \in S} \sum_{w \in T} c(v \rightarrow w) && \text{since } f(u \rightarrow v) \leq c(v \rightarrow w) \\
 &= \|S, T\| && \text{by definition}
 \end{aligned}$$

Our derivation actually implies the following stronger observation:  $|f| = \|S, T\|$  if and only if  $f$  **saturnates every edge from  $S$  to  $T$  and avoids every edge from  $T$  to  $S$** . Moreover, if we have a flow  $f$  and a cut  $(S, T)$  that satisfies this equality condition,  $f$  must be a maximum flow, and  $(S, T)$  must be a minimum cut.

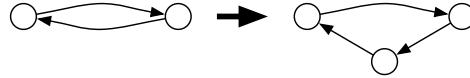
## 21.3 The Max-Flow Min-Cut Theorem

Surprisingly, for any weighted directed graph, there is always a flow  $f$  and a cut  $(S, T)$  that satisfy the equality condition. This is the famous *max-flow min-cut theorem*:

**The value of the maximum flow is equal to the capacity of the minimum cut.**

The rest of this section gives a proof of this theorem; we will eventually turn this proof into an algorithm.

Fix a graph  $G$ , vertices  $s$  and  $t$ , and a capacity function  $c: E \rightarrow \mathbb{R}_{\geq 0}$ . The proof will be easier if we assume that the capacity function is **reduced**: For any vertices  $u$  and  $v$ , either  $c(u \rightarrow v) = 0$  or  $c(v \rightarrow u) = 0$ , or equivalently, if an edge appears in  $G$ , then its reversal does not. This assumption is easy to enforce. Whenever an edge  $u \rightarrow v$  and its reversal  $v \rightarrow u$  are both the graph, replace the edge  $u \rightarrow v$  with a path  $u \rightarrow x \rightarrow v$  of length two, where  $x$  is a new vertex and  $c(u \rightarrow x) = c(x \rightarrow v) = c(u \rightarrow v)$ . The modified graph has the same maximum flow value and minimum cut capacity as the original graph.

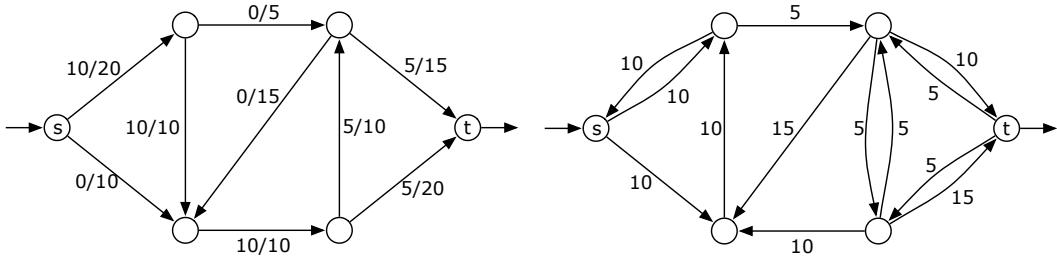


Enforcing the one-direction assumption.

Let  $f$  be a feasible flow. We define a new capacity function  $c_f: V \times V \rightarrow \mathbb{R}$ , called the **residual capacity**, as follows:

$$c_f(u \rightarrow v) = \begin{cases} c(u \rightarrow v) - f(u \rightarrow v) & \text{if } u \rightarrow v \in E \\ f(v \rightarrow u) & \text{if } v \rightarrow u \in E \\ 0 & \text{otherwise} \end{cases}$$

Since  $f \geq 0$  and  $f \leq c$ , the residual capacities are always non-negative. It is possible to have  $c_f(u \rightarrow v) > 0$  even if  $u \rightarrow v$  is not an edge in the original graph  $G$ . Thus, we define the **residual graph**  $G_f = (V, E_f)$ , where  $E_f$  is the set of edges whose residual capacity is positive. Notice that the residual capacities are *not* necessarily reduced; it is quite possible to have both  $c_f(u \rightarrow v) > 0$  and  $c_f(v \rightarrow u) > 0$ .



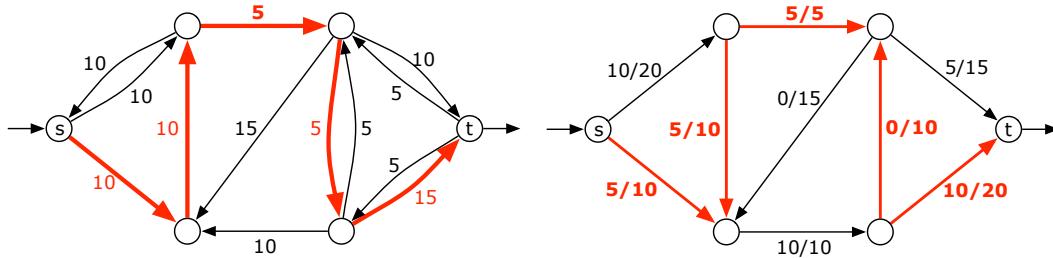
A flow  $f$  in a weighted graph  $G$  and the corresponding residual graph  $G_f$ .

Suppose there is no path from the source  $s$  to the target  $t$  in the residual graph  $G_f$ . Let  $S$  be the set of vertices that are reachable from  $s$  in  $G_f$ , and let  $T = V \setminus S$ . The partition  $(S, T)$  is clearly an  $(s, t)$ -cut. For every vertex  $u \in S$  and  $v \in T$ , we have

$$c_f(u \rightarrow v) = (c(u \rightarrow v) - f(u \rightarrow v)) + f(v \rightarrow u) = 0,$$

which implies that  $c(u \rightarrow v) - f(u \rightarrow v) = 0$  and  $f(v \rightarrow u) = 0$ . In other words, our flow  $f$  saturates every edge from  $S$  to  $T$  and avoids every edge from  $T$  to  $S$ . It follows that  $|f| = \|S, T\|$ . Moreover,  $f$  is a maximum flow and  $(S, T)$  is a minimum cut.

On the other hand, suppose there is a path  $s = v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_r = t$  in  $G_f$ . We refer to  $v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_r$  as an **augmenting path**. Let  $F = \min_i c_f(v_i \rightarrow v_{i+1})$  denote the maximum amount of flow that we can



An augmenting path in  $G_f$  with value  $F = 5$  and the augmented flow  $f'$ .

push through the augmenting path in  $G_f$ . We define a new flow function  $f': E \rightarrow \mathbb{R}$  as follows:

$$f'(u \rightarrow v) = \begin{cases} f(u \rightarrow v) + F & \text{if } u \rightarrow v \text{ is in the augmenting path} \\ f(u \rightarrow v) - F & \text{if } v \rightarrow u \text{ is in the augmenting path} \\ f(u \rightarrow v) & \text{otherwise} \end{cases}$$

To prove that the flow  $f'$  is feasible with respect to the original capacities  $c$ , we need to verify that  $f' \geq 0$  and  $f' \leq c$ . Consider an edge  $u \rightarrow v$  in  $G$ . If  $u \rightarrow v$  is in the augmenting path, then  $f'(u \rightarrow v) > f(u \rightarrow v) \geq 0$  and

$$\begin{aligned} f'(u \rightarrow v) &= f(u \rightarrow v) + F && \text{by definition of } f' \\ &\leq f(u \rightarrow v) + c_f(u \rightarrow v) && \text{by definition of } F \\ &= f(u \rightarrow v) + c(u \rightarrow v) - f(u \rightarrow v) && \text{by definition of } c_f \\ &= c(u \rightarrow v) && \text{Duh.} \end{aligned}$$

On the other hand, if the reversal  $v \rightarrow u$  is in the augmenting path, then  $f'(u \rightarrow v) < f(u \rightarrow v) \leq c(u \rightarrow v)$ , which implies that

$$\begin{aligned} f'(u \rightarrow v) &= f(u \rightarrow v) - F && \text{by definition of } f' \\ &\geq f(u \rightarrow v) - c_f(v \rightarrow u) && \text{by definition of } F \\ &= f(u \rightarrow v) - f(u \rightarrow v) && \text{by definition of } c_f \\ &= 0 && \text{Duh.} \end{aligned}$$

Finally, we observe that (without loss of generality) only the first edge in the augmenting path leaves  $s$ , so  $|f'| = |f| + F > 0$ . In other words,  $f$  is *not* a maximum flow.

This completes the proof!

## Exercises

- Let  $(S, T)$  and  $(S', T')$  be minimum  $(s, t)$ -cuts in some flow network  $G$ . Prove that  $(S \cap S', T \cup T')$  and  $(S \cup S', T \cap T')$  are also minimum  $(s, t)$ -cuts in  $G$ .
- Suppose  $(S, T)$  is the *unique* minimum  $(s, t)$ -cut in some flow network. Prove that  $(S, T)$  is also a minimum  $(x, y)$ -cut for all vertices  $x \in S$  and  $y \in T$ .

3. Cuts are sometimes defined as subsets of the edges of the graph, instead of as partitions of its vertices. In this problem, you will prove that these two definitions are *almost* equivalent.

We say that a subset  $X$  of (directed) edges *separates*  $s$  and  $t$  if every directed path from  $s$  to  $t$  contains at least one (directed) edge in  $X$ . For any subset  $S$  of vertices, let  $\delta S$  denote the set of directed edges leaving  $S$ ; that is,  $\delta S := \{u \rightarrow v \mid u \in S, v \notin S\}$ .

- (a) Prove that if  $(S, T)$  is an  $(s, t)$ -cut, then  $\delta S$  separates  $s$  and  $t$ .
  - (b) Let  $X$  be an arbitrary subset of edges that separates  $s$  and  $t$ . Prove that there is an  $(s, t)$ -cut  $(S, T)$  such that  $\delta S \subseteq X$ .
  - (c) Let  $X$  be a *minimal* subset of edges that separates  $s$  and  $t$ . Prove that there is an  $(s, t)$ -cut  $(S, T)$  such that  $\delta S = X$ .
4. A flow  $f$  is **acyclic** if the subgraph of directed edges with positive flow contains no directed cycles.
- (a) Prove that for any flow  $f$ , there is an acyclic flow with the same value as  $f$ . (In particular, this implies that some maximum flow is acyclic.)
  - (b) A *path flow* assigns positive values only to the edges of one simple directed path from  $s$  to  $t$ . Prove that every acyclic flow can be written as the sum of  $O(E)$  path flows.
  - (c) Describe a flow in a directed graph that *cannot* be written as the sum of path flows.
  - (d) A *cycle flow* assigns positive values only to the edges of one simple directed cycle. Prove that every flow can be written as the sum of  $O(E)$  path flows and cycle flows.
  - (e) Prove that every flow with value 0 can be written as the sum of  $O(E)$  cycle flows. (Zero-value flows are also called *circulations*.)
5. Suppose instead of capacities, we consider networks where each edge  $u \rightarrow v$  has a non-negative **demand**  $d(u \rightarrow v)$ . Now an  $(s, t)$ -flow  $f$  is *feasible* if and only if  $f(u \rightarrow v) \geq d(u \rightarrow v)$  for every edge  $u \rightarrow v$ . (Feasible flow values can now be arbitrarily large.) A natural problem in this setting is to find a feasible  $(s, t)$ -flow of *minimum* value.
- (a) Describe an efficient algorithm to compute a feasible  $(s, t)$ -flow, given the graph, the demand function, and the vertices  $s$  and  $t$  as input. [*Hint: Find a flow that is non-zero everywhere, and then scale it up to make it feasible.*]
  - (b) Suppose you have access to a subroutine `MAXFLOW` that computes *maximum* flows in networks with edge capacities. Describe an efficient algorithm to compute a *minimum* flow in a given network with edge demands; your algorithm should call `MAXFLOW` exactly once.
  - (c) State and prove an analogue of the max-flow min-cut theorem for this setting. (Do minimum flows correspond to maximum cuts?)

*A process cannot be understood by stopping it. Understanding must move with the flow of the process, must join it and flow with it.*

— The First Law of Mentat, in Frank Herbert's *Dune* (1965)

*There's a difference between knowing the path and walking the path.*

— Morpheus [Laurence Fishburne], *The Matrix* (1999)

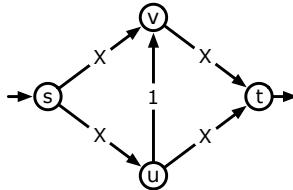
## 22 Max-Flow Algorithms

### 22.1 Ford-Fulkerson augmenting path

Ford and Fulkerson's proof of the Maxflow Mincut theorem, described in the previous lecture note, translates immediately to an algorithm to compute maximum flows: Starting with the zero flow, repeatedly augment the flow along **any** path  $s \rightsquigarrow t$  in the residual graph, until there is no such path.

If every edge capacity is an integer, then every augmentation step increases the value of the flow by a positive integer. Thus, the algorithm halts after  $|f^*|$  iterations, where  $f^*$  is the actual maximum flow. Each iteration requires  $O(E)$  time, to create the residual graph  $G_f$  and perform a whatever-first-search to find an augmenting path. Thus the Ford-Fulkerson algorithm runs in  $O(E|f^*|)$  time in the worst case.

The following example shows that this running time analysis is essentially tight. Consider the 4-node network illustrated below, where  $X$  is some large integer. The maximum flow in this network is clearly  $2X$ . Moreover, the Ford-Fulkerson could alternate between pushing 1 unit of flow along the augmenting path  $s \rightarrow u \rightarrow v \rightarrow t$  and then pushing 1 unit of flow along the augmenting path  $s \rightarrow v \rightarrow u \rightarrow t$ , leading to a running time of  $\Theta(X) = \Omega(|f^*|)$ .

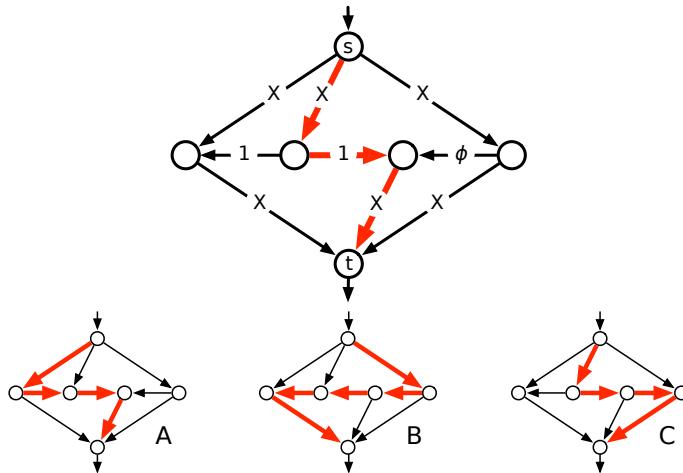


A bad example for the Ford-Fulkerson algorithm.

Ford and Fulkerson's algorithm works quite well in many practical situations, or in settings where the maximum flow value  $|f^*|$  is small, but without further constraints on the augmenting paths, this is *not* an efficient algorithm in general. The example network above can be described using only  $O(\log X)$  bits; thus, the running time of Ford-Fulkerson is *exponential* in the input size.

If we multiply all the capacities by the same (positive) constant, the maximum flow increases everywhere by the same constant factor. It follows that if all the edge capacities are *rational*, then the Ford-Fulkerson algorithm eventually halts, although still in exponential time. However, if we allow irrational capacities, the algorithm can actually loop forever, always finding smaller and smaller augmenting paths! Worse yet, this infinite sequence of augmentations may not even converge to the maximum flow! Perhaps the simplest example of this effect was discovered by Uri Zwick.

Consider the six-node network shown on the next page. Six of the nine edges have some large integer capacity  $X$ , two have capacity 1, and one has capacity  $\phi = (\sqrt{5} - 1)/2 \approx 0.618034$ , chosen so that  $1 - \phi = \phi^2$ . To prove that the Ford-Fulkerson algorithm can get stuck, we can watch the residual capacities of the three horizontal edges as the algorithm progresses. (The residual capacities of the other six edges will always be at least  $X - 3$ .)



Uri Zwick's non-terminating flow example, and three augmenting paths.

The Ford-Fulkerson algorithm starts by choosing the central augmenting path, shown in the large figure above. The three horizontal edges, in order from left to right, now have residual capacities 1, 0, and  $\phi$ . Suppose inductively that the horizontal residual capacities are  $\phi^{k-1}$ , 0,  $\phi^k$  for some non-negative integer  $k$ .

1. Augment along  $B$ , adding  $\phi^k$  to the flow; the residual capacities are now  $\phi^{k+1}$ ,  $\phi^k$ , 0.
2. Augment along  $C$ , adding  $\phi^k$  to the flow; the residual capacities are now  $\phi^{k+1}$ , 0,  $\phi^k$ .
3. Augment along  $B$ , adding  $\phi^{k+1}$  to the flow; the residual capacities are now 0,  $\phi^{k+1}$ ,  $\phi^{k+2}$ .
4. Augment along  $A$ , adding  $\phi^{k+1}$  to the flow; the residual capacities are now  $\phi^{k+1}$ , 0,  $\phi^{k+2}$ .

Thus, after  $4n + 1$  augmentation steps, the residual capacities are  $\phi^{2n-2}$ , 0,  $\phi^{2n-1}$ . As the number of augmentation steps grows to infinity, the value of the flow converges to

$$1 + 2 \sum_{i=1}^{\infty} \phi^i = 1 + \frac{2}{1 - \phi} = 4 + \sqrt{5} < 7,$$

even though the maximum flow value is clearly  $2X + 1$ .

Picky students might wonder at this point why we care about irrational capacities; after all, computers can't represent anything but (small) integers or (dyadic) rationals exactly. Good question! One reason is that the integer restriction is literally *artificial*; it's an *artifact* of actual computational hardware<sup>1</sup>, not an inherent feature of the abstract mathematical problem. Another reason, which is probably more convincing to most practical computer scientists, is that the behavior of the algorithm with irrational inputs tells us something about its worst-case behavior *in practice* given floating-point capacities—terrible! Even with very reasonable capacities, a careless implementation of Ford-Fulkerson could enter an infinite loop simply because of round-off error!

## 22.2 Edmonds-Karp: Fat Pipes

The Ford-Fulkerson algorithm does not specify which alternating path to use if there is more than one. In 1972, Jack Edmonds and Richard Karp analyzed two natural heuristics for choosing the path. The first is essentially a greedy algorithm:

<sup>1</sup> ...or perhaps the laws of physics. Yeah, right. Whatever. Like reality actually matters in this class.

Choose the augmenting path with largest bottleneck value.

It's a fairly easy to show that the maximum-bottleneck  $(s, t)$ -path in a directed graph can be computed in  $O(E \log V)$  time using a variant of Jarník's minimum-spanning-tree algorithm, or of Dijkstra's shortest path algorithm. Simply grow a directed spanning tree  $T$ , rooted at  $s$ . Repeatedly find the highest-capacity edge leaving  $T$  and add it to  $T$ , until  $T$  contains a path from  $s$  to  $t$ . Alternately, one could emulate Kruskal's algorithm—insert edges one at a time in decreasing capacity order until there is a path from  $s$  to  $t$ —although this is less efficient.

We can now analyze the algorithm in terms of the value of the maximum flow  $f^*$ . Let  $f$  be any flow in  $G$ , and let  $f'$  be the maximum flow in the current residual graph  $G_f$ . (At the beginning of the algorithm,  $G_f = G$  and  $f' = f^*$ .) Let  $e$  be the bottleneck edge in the next augmenting path. Let  $S$  be the set of vertices reachable from  $s$  through edges in  $G$  with capacity greater than  $c(e)$  and let  $T = V \setminus S$ . By construction,  $T$  is non-empty, and every edge from  $S$  to  $T$  has capacity at most  $c(e)$ . Thus, the capacity of the cut  $(S, T)$  is at most  $c(e) \cdot |E|$ . On the other hand, the maxflow-mincut theorem implies that  $\|S, T\| \geq |f|$ . We conclude that  $c(e) \geq |f|/|E|$ .

The preceding argument implies that augmenting  $f$  along the maximum-bottleneck path in  $G_f$  multiplies the maximum flow value in  $G_f$  by a factor of at most  $1 - 1/E$ . In other words, the residual flow decays exponentially with the number of iterations. After  $E \cdot \ln|f^*|$  iterations, the maximum flow value in  $G_f$  is at most

$$|f^*| \cdot (1 - 1/E)^{E \cdot \ln|f^*|} < |f^*| e^{-\ln|f^*|} = 1.$$

(That's Euler's constant  $e$ , not the edge  $e$ . Sorry.) In particular, if all the capacities are integers, then after  $E \cdot \ln|f^*|$  iterations, the maximum capacity of the residual graph is zero and  $f$  is a maximum flow.

We conclude that for graphs with integer capacities, the Edmonds-Karp ‘fat pipe’ algorithm runs in  $O(E^2 \log E \log |f^*|)$  time, which is actually a polynomial function of the input size.

### 22.3 Dinitz/Edmonds-Karp: Short Pipes

The second Edmonds-Karp heuristic was actually proposed by Ford and Fulkerson in their original max-flow paper, and first analyzed by the Russian mathematician Dinitz (sometimes transliterated Dinic) in 1970. Edmonds and Karp published their independent and slightly weaker analysis in 1972. So naturally, almost everyone refers to this algorithm as ‘Edmonds-Karp’.<sup>2</sup>

Choose the augmenting path with fewest edges.

The correct path can be found in  $O(E)$  time by running breadth-first search in the residual graph. More surprisingly, the algorithm halts after a polynomial number of iterations, independent of the actual edge capacities!

The proof of this upper bound relies on two observations about the evolution of the residual graph. Let  $f_i$  be the current flow after  $i$  augmentation steps, let  $G_i$  be the corresponding residual graph. In particular,  $f_0$  is zero everywhere and  $G_0 = G$ . For each vertex  $v$ , let  $\text{level}_i(v)$  denote the unweighted shortest path distance from  $s$  to  $v$  in  $G_i$ , or equivalently, the *level* of  $v$  in a breadth-first search tree of  $G_i$  rooted at  $s$ .

Our first observation is that these levels can only increase over time.

<sup>2</sup>To be fair, Edmonds and Karp discovered their algorithm a few years before publication—getting ideas into print takes time, especially in the early 1970s—which is why some authors believe they deserve priority. I don't buy it; Dinitz also presumably discovered his algorithm a few years before its publication. (In Soviet Union, result publish you.) On the gripping hand, Dinitz's paper also described an improvement to the algorithm presented here that runs in  $O(V^2E)$  time instead of  $O(VE^2)$ , so maybe that ought to be called Dinitz's algorithm.

**Lemma 1.**  $\text{level}_{i+1}(v) \geq \text{level}_i(v)$  for all vertices  $v$  and integers  $i$ .

**Proof:** The claim is trivial for  $v = s$ , since  $\text{level}_i(s) = 0$  for all  $i$ . Choose an arbitrary vertex  $v \neq s$ , and let  $s \rightarrow \dots \rightarrow u \rightarrow v$  be a shortest path from  $s$  to  $v$  in  $G_{i+1}$ . (If there is no such path, then  $\text{level}_{i+1}(v) = \infty$ , and we're done.) Because this is a shortest path, we have  $\text{level}_{i+1}(v) = \text{level}_{i+1}(u) + 1$ , and the inductive hypothesis implies that  $\text{level}_{i+1}(u) \geq \text{level}_i(u)$ .

We now have two cases to consider. If  $u \rightarrow v$  is an edge in  $G_i$ , then  $\text{level}_i(v) \leq \text{level}_i(u) + 1$ , because the levels are defined by breadth-first traversal.

On the other hand, if  $u \rightarrow v$  is not an edge in  $G_i$ , then  $v \rightarrow u$  must be an edge in the  $i$ th augmenting path. Thus,  $v \rightarrow u$  must lie on the shortest path from  $s$  to  $t$  in  $G_i$ , which implies that  $\text{level}_i(v) = \text{level}_i(u) - 1 \leq \text{level}_i(u) + 1$ .

In both cases, we have  $\text{level}_{i+1}(v) = \text{level}_{i+1}(u) + 1 \geq \text{level}_i(u) + 1 \geq \text{level}_i(v)$ .  $\square$

Whenever we augment the flow, the bottleneck edge in the augmenting path disappears from the residual graph, and some other edge in the *reversal* of the augmenting path may (re-)appear. Our second observation is that an edge cannot appear or disappear too many times.

**Lemma 2.** During the execution of the Dinitz/Edmonds-Karp algorithm, any edge  $u \rightarrow v$  disappears from the residual graph  $G_f$  at most  $V/2$  times.

**Proof:** Suppose  $u \rightarrow v$  is in two residual graphs  $G_i$  and  $G_{j+1}$ , but not in any of the intermediate residual graphs  $G_{i+1}, \dots, G_j$ , for some  $i < j$ . Then  $u \rightarrow v$  must be in the  $i$ th augmenting path, so  $\text{level}_i(v) = \text{level}_i(u) + 1$ , and  $v \rightarrow u$  must be on the  $j$ th augmenting path, so  $\text{level}_j(v) = \text{level}_j(u) - 1$ . By the previous lemma, we have

$$\text{level}_j(u) = \text{level}_j(v) + 1 \geq \text{level}_i(v) + 1 = \text{level}_i(u) + 2.$$

In other words, the distance from  $s$  to  $u$  increased by at least 2 between the disappearance and reappearance of  $u \rightarrow v$ . Since every level is either less than  $V$  or infinite, the number of disappearances is at most  $V/2$ .  $\square$

Now we can derive an upper bound on the number of iterations. Since each edge can disappear at most  $V/2$  times, there are at most  $EV/2$  edge disappearances overall. But at least one edge disappears on each iteration, so the algorithm must halt after at most  $EV/2$  iterations. Finally, since each iteration requires  $O(E)$  time, Dinitz' algorithm runs in  **$O(VE^2)$  time** overall.

## 22.4 Push-Relabel

To be written!

## 22.5 Further Progress

This is nowhere near the end of the story for maximum-flow algorithms. Decades of further research have led to a number of even faster algorithms, some of which are summarized in the table below.<sup>3</sup> All of the algorithms listed below compute a maximum flow in several iterations. Each algorithm has two variants: a simpler version that performs each iteration by brute force, and a faster variant that uses sophisticated data structures to maintain a spanning tree of the flow network, so that each iteration can be performed (and the spanning tree updated) in logarithmic time. There is no reason to believe that the best algorithms known so far are optimal; indeed, maximum flows are still a very active area of research.

<sup>3</sup>To keep the table short, I have deliberately omitted algorithms whose running time depends on the maximum capacity, the sum of the capacities, or the maximum flow value. Even with this restriction, the table is incomplete!

Technique	Direct	With dynamic trees	Sources
Blocking flow	$O(V^3)$	$O(VE \log V)$	[Dinitz; Sleator and Tarjan]
Network simplex	$O(V^2E)$	$O(VE \log V)$	[Dantzig; Goldfarb and Hao; Goldberg, Grigoriadis, and Tarjan]
Push-relabel (generic)	$O(V^2E)$	—	[Goldberg and Tarjan]
Push-relabel (FIFO)	$O(V^3)$	$O(V^2 \log(V^2/E))$	[Goldberg and Tarjan]
Push-relabel (highest label)	$O(V^2\sqrt{E})$	—	[Cheriyan and Maheshwari; Tunçel]
Pseudoflow	$O(V^2E)$	$O(VE \log V)$	[Hochbaum]

Several purely combinatorial maximum-flow algorithms and their running times.

## Exercises

1. A new assistant professor, teaching maximum flows for the first time, suggests the following greedy modification to the generic Ford-Fulkerson augmenting path algorithm. Instead of maintaining a residual graph, just reduce the capacity of edges along the augmenting path! In particular, whenever we saturate an edge, just remove it from the graph.

```

GREEDYFLOW( $G, c, s, t$ ):
  for every edge  $e$  in  $G$ 
     $f(e) \leftarrow 0$ 
  while there is a path from  $s$  to  $t$ 
     $\pi \leftarrow$  an arbitrary path from  $s$  to  $t$ 
     $F \leftarrow$  minimum capacity of any edge in  $\pi$ 
    for every edge  $e$  in  $\pi$ 
       $f(e) \leftarrow f(e) + F$ 
      if  $c(e) = F$ 
        remove  $e$  from  $G$ 
      else
         $c(e) \leftarrow c(e) - F$ 
  return  $f$ 

```

- (a) Show that this algorithm does *not* always compute a maximum flow.
- (b) Prove that for any flow network, if the Greedy Path Fairy tells you precisely which path  $\pi$  to use at each iteration, then GREEDYFLOW does compute a maximum flow. (Sadly, the Greedy Path Fairy does not actually exist.)
- 2. Describe and analyze an algorithm to find the maximum-bottleneck path from  $s$  to  $t$  in a flow network  $G$  in  $O(E \log V)$  time.
- ★3. Describe a directed graph with irrational edge capacities, such that the Edmonds-Karp ‘fat pipe’ heuristic does not halt.
- 4. Describe an efficient algorithm to check whether a given flow network contains a *unique* maximum flow.
- 5. For any flow network  $G$  and any vertices  $u$  and  $v$ , let  $bottleneck_G(u, v)$  denote the maximum, over all paths  $\pi$  in  $G$  from  $u$  to  $v$ , of the minimum-capacity edge along  $\pi$ . Describe an algorithm

to construct a spanning tree  $T$  of  $G$  such that  $bottleneck_T(u, v) = bottleneck_G(u, v)$ . (Edges in  $T$  inherit their capacities from  $G$ .)

One way to think about this problem is to imagine the vertices of the graph as islands, and the edges as bridges. Each bridge has a maximum weight it can support. If a truck is carrying stuff from  $u$  to  $v$ , how much can the truck carry? We don't care what route the truck takes; the point is that the smallest-weight edge on the route will determine the load.

6. We can speed up the Edmonds-Karp 'fat pipe' heuristic, at least for integer capacities, by relaxing our requirements for the next augmenting path. Instead of finding the augmenting path with maximum bottleneck capacity, we find a path whose bottleneck capacity is at least half of maximum, using the following *capacity scaling* algorithm.

The algorithm maintains a bottleneck threshold  $\Delta$ ; initially,  $\Delta$  is the maximum capacity among all edges in the graph. In each *phase*, the algorithm augments along paths from  $s$  to  $t$  in which every edge has residual capacity at least  $\Delta$ . When there is no such path, the phase ends, we set  $\Delta \leftarrow \Delta/2$ , and the next phase begins.

- (a) How many phases will the algorithm execute in the worst case, if the edge capacities are integers?
- (b) Let  $f$  be the flow at the end of a phase for a particular value of  $\Delta$ . Let  $S$  be the nodes that are reachable from  $s$  in the residual graph  $G_f$  using only edges with residual capacity at least  $\Delta$ , and let  $T = V \setminus S$ . Prove that the capacity (with respect to  $G$ 's original edge capacities) of the cut  $(S, T)$  is at most  $|f| + E \cdot \Delta$ .
- (c) Prove that in each phase of the scaling algorithm, there are at most  $2E$  augmentations.
- (d) What is the overall running time of the scaling algorithm, assuming all the edge capacities are integers?

*For a long time it puzzled me how something so expensive, so leading edge, could be so useless, and then it occurred to me that a computer is a stupid machine with the ability to do incredibly smart things, while computer programmers are smart people with the ability to do incredibly stupid things. They are, in short, a perfect match.*

— Bill Bryson, *Notes from a Big Country* (1999)

## 23 Applications of Maximum Flow

### 23.1 Edge-Disjoint Paths

One of the easiest applications of maximum flows is computing the maximum number of edge-disjoint paths between two specified vertices  $s$  and  $t$  in a directed graph  $G$  using maximum flows. A set of paths in  $G$  is *edge-disjoint* if each edge in  $G$  appears in at most one of the paths; several edge-disjoint paths may pass through the same vertex, however.

If we give each edge capacity 1, then the maxflow from  $s$  to  $t$  assigns a flow of either 0 or 1 to every edge. Since any vertex of  $G$  lies on at most two saturated edges (one in and one out, or none at all), the subgraph  $S$  of saturated edges is the union of several edge-disjoint paths and cycles. Moreover, the number of paths is exactly equal to the value of the flow. Extracting the actual paths from  $S$  is easy—just follow any directed path in  $S$  from  $s$  to  $t$ , remove that path from  $S$ , and recurse.

Conversely, we can transform any collection of  $k$  edge-disjoint paths into a flow by pushing one unit of flow along each path from  $s$  to  $t$ ; the value of the resulting flow is exactly  $k$ . It follows that the maxflow algorithm actually computes the largest possible set of edge-disjoint paths. The overall running time is  $O(VE)$ , just like for maximum bipartite matchings.

The same algorithm can also be used to find edge-disjoint paths in *undirected* graphs. We simply replace every undirected edge in  $G$  with a pair of directed edges, each with unit capacity, and compute a maximum flow from  $s$  to  $t$  in the resulting directed graph  $G'$  using the Ford-Fulkerson algorithm. For any edge  $uv$  in  $G$ , if our max flow saturates both directed edges  $u \rightarrow v$  and  $v \rightarrow u$  in  $G'$ , we can remove *both* edges from the flow without changing its value. Thus, without loss of generality, the maximum flow assigns a direction to every saturated edge, and we can extract the edge-disjoint paths by searching the graph of directed saturated edges.

### 23.2 Vertex Capacities and Vertex-Disjoint Paths

Suppose we have capacities on the vertices as well as the edges. Here, in addition to our other constraints, we require that for any vertex  $v$  other than  $s$  and  $t$ , the total flow into  $v$  (and therefore the total flow out of  $v$ ) is at most some non-negative value  $c(v)$ . How can we compute a maximum flow with these new constraints?

One possibility is to modify our existing algorithms to take these vertex capacities into account. Given a flow  $f$ , we can define the *residual capacity* of a vertex  $v$  to be its original capacity minus the total flow into  $v$ :

$$c_f(v) = c(v) - \sum_u f(u \rightarrow v).$$

Since we cannot send any more flow into a vertex with residual capacity 0 we remove from the residual graph  $G_f$  every edge  $u \rightarrow v$  that appears in  $G$  whose head vertex  $v$  is saturated. Otherwise, the augmenting-path algorithm is unchanged.

But an even simpler method is to transform the input into a traditional flow network, with only edge capacities. Specifically, we replace every vertex  $v$  with two vertices  $v_{\text{in}}$  and  $v_{\text{out}}$ , connected by an edge  $v_{\text{in}} \rightarrow v_{\text{out}}$  with capacity  $c(v)$ , and then replace every directed edge  $u \rightarrow v$  with the edge  $u_{\text{out}} \rightarrow v_{\text{in}}$  (keeping the same capacity). Finally, we compute the maximum flow from  $s_{\text{out}}$  to  $t_{\text{in}}$  in this modified flow network.

It is now easy to compute the maximum number of vertex-disjoint paths from  $s$  to  $t$  in any directed graph. Simply give every vertex capacity 1, and compute a maximum flow!

### 23.3 Maximum Matchings in Bipartite Graphs

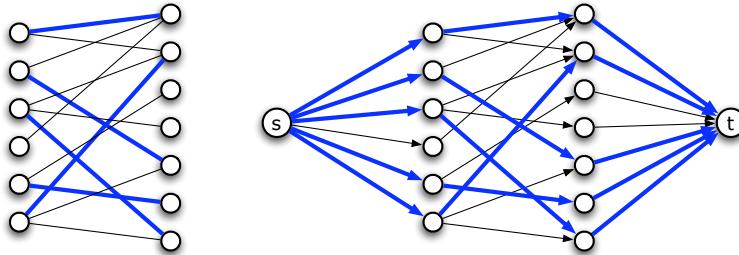
Another natural application of maximum flows is finding large *matchings* in bipartite graphs. A matching is a subgraph in which every vertex has degree at most one, or equivalently, a collection of edges such that no two share a vertex. The problem is to find the matching with the maximum number of edges in a given bipartite graph.

We can solve this problem by reducing it to a maximum flow problem as follows. Let  $G$  be the given bipartite graph with vertex set  $U \cup W$ , such that every edge joins a vertex in  $U$  to a vertex in  $W$ . We create a new *directed* graph  $G'$  by (1) orienting each edge from  $U$  to  $W$ , (2) adding two new vertices  $s$  and  $t$ , (3) adding edges from  $s$  to every vertex in  $U$ , and (4) adding edges from each vertex in  $W$  to  $t$ . Finally, we assign every edge in  $G'$  a capacity of 1.

Any matching  $M$  in  $G$  can be transformed into a flow  $f_M$  in  $G'$  as follows: For each edge  $uw$  in  $M$ , push one unit of flow along the path  $s \rightarrow u \rightarrow w \rightarrow t$ . These paths are disjoint except at  $s$  and  $t$ , so the resulting flow satisfies the capacity constraints. Moreover, the value of the resulting flow is equal to the number of edges in  $M$ .

Conversely, consider any  $(s, t)$ -flow  $f$  in  $G'$  computed using the Ford-Fulkerson augmenting path algorithm. Because the edge capacities are integers, the Ford-Fulkerson algorithm assigns an integer flow to every edge. (This is easy to verify by induction, hint, hint.) Moreover, since each edge has *unit* capacity, the computed flow either saturates ( $f(e) = 1$ ) or avoids ( $f(e) = 0$ ) every edge in  $G'$ . Finally, since at most one unit of flow can enter any vertex in  $U$  or leave any vertex in  $W$ , the saturated edges from  $U$  to  $W$  form a matching in  $G$ . The size of this matching is exactly  $|f|$ .

Thus, the size of the maximum matching in  $G$  is equal to the value of the maximum flow in  $G'$ , and provided we compute the maxflow using augmenting paths, we can convert the actual maxflow into a maximum matching. The maximum flow has value at most  $\min\{|U|, |W|\} = O(V)$ , so the Ford-Fulkerson algorithm runs in  $O(VE)$  time.



A maximum matching in a bipartite graph  $G$ , and the corresponding maximum flow in  $G'$ .

## 23.4 Binary Assignment Problems

Maximum-cardinality matchings are a special case of a general family of so-called *assignment* problems.<sup>1</sup> An unweighted *binary* assignment problem involves two disjoint finite sets  $X$  and  $Y$ , which typically represent two different kinds of resources, such as web pages and servers, jobs and machines, rows and columns of a matrix, hospitals and interns, or customers and pints of ice cream. Our task is to choose the largest possible collection of pairs  $(x, y)$  as possible, where  $x \in X$  and  $y \in Y$ , subject to several constraints of the following form:

- Each element  $x \in X$  can appear in at most  $c(x)$  pairs.
- Each element  $y \in Y$  can appear in at most  $c(y)$  pairs.
- Each pair  $(x, y) \in X \times Y$  can appear in the output at most  $c(x, y)$  times.

Each upper bound  $c(x)$ ,  $c(y)$ , and  $c(x, y)$  is either a (typically small) non-negative integer or  $\infty$ . Intuitively, we create each pair in our output by *assigning* an element of  $X$  to an element of  $Y$ .

The maximum-matching problem is a special case, where  $c(z) = 1$  for all  $z \in X \cup Y$ , and each  $c(x, y)$  is either 0 or 1, depending on whether the pair  $xy$  defines an edge in the underlying bipartite graph.

Here is a slightly more interesting example. A nearby school, famous for its onerous administrative hurdles, decides to organize a dance. Every pair of students (one boy, one girl) who wants to dance must register in advance. School regulations limit each boy-girl pair to at most three dances together, and limits each student to at most ten dances overall. How can we maximize the number of dances? This is a binary assignment problem for the set  $X$  of girls and the set  $Y$  of boys, where for each girl  $x$  and boy  $y$ , we have  $c(x) = c(y) = 10$  and either  $c(x, y) = 3$  (if  $x$  and  $y$  registered to dance) or  $c(x, y) = 0$  (if they didn't register).

Every binary assignment problem can be reduced to a standard maximum flow problem as follows. We construct a flow network  $G = (V, E)$  with vertices  $X \cup Y \cup \{s, t\}$  and the following edges:

- an edge  $s \rightarrow x$  with capacity  $c(x)$  for each  $x \in X$ ,
- an edge  $y \rightarrow t$  with capacity  $c(y)$  for each  $y \in Y$ .
- an edge  $x \rightarrow y$  with capacity  $c(x, y)$  for each  $x \in X$  and  $y \in Y$ , and

Because all the edges have integer capacities, the Ford-Fulkerson algorithm constructs an integer maximum flow  $f^*$ . This flow can be decomposed into the sum of  $|f^*|$  paths of the form  $s \rightarrow x \rightarrow y \rightarrow t$  for some  $x \in X$  and  $y \in Y$ . For each such path, we report the pair  $(x, y)$ . (Equivalently, the pair  $(x, y)$  appears in our output collection  $f(x \rightarrow y)$  times.) It is easy to verify (hint, hint) that this collection of pairs satisfies all the necessary constraints. Conversely, any legal collection of  $r$  pairs can be transformed into a feasible integer flow with value  $r$  in  $G$ . Thus, the largest legal collection of pairs corresponds to a maximum flow in  $G$ . So our algorithm is correct.

The value of the maximum flow clearly is at most  $CV$ , where  $C$  is an upper bound on all the vertex and edge capacities, and  $n$  is the number of vertices, so the off-the-shelf Ford-Fulkerson algorithm runs in  $O(E|f^*|) = O(CVE)$  time. If  $C = O(1)$ , this is more efficient than any specialized flow algorithm.

## 23.5 Baseball Elimination

Every year millions of baseball fans eagerly watch their favorite team, hoping they will win a spot in the playoffs, and ultimately the World Series. Sadly, most teams are “mathematically eliminated” days or

---

<sup>1</sup>Most authors refer to finding a maximum-weight matching in a bipartite graph as *the* assignment problem.

even weeks before the regular season ends. Often, it is easy to spot when a team is eliminated—they can't win enough games to catch up to the current leader in their division. But sometimes the situation is more subtle.

For example, here are the actual standings from the American League East on August 30, 1996.

Team	Won-Lost	Left	NYY	BAL	BOS	TOR	DET
New York Yankees	75–59	28		3	8	7	3
Baltimore Orioles	71–63	28	3		2	7	4
Boston Red Sox	69–66	27	8	2		0	0
Toronto Blue Jays	63–72	27	7	7	0		0
Detroit Lions	49–86	27	3	4	0	0	

Detroit is clearly behind, but some die-hard Lions fans may hold out hope that their team can still win. After all, if Detroit wins all 27 of their remaining games, they will end the season with 76 wins, more than any other team has now. So as long as every other team loses every game... but that's not possible, because some of those other teams still have to play each other. Here is one complete argument:<sup>2</sup>

*By winning all of their remaining games, Detroit can finish the season with a record of 76 and 86. If the Yankees win just 2 more games, then they will finish the season with a 77 and 85 record which would put them ahead of Detroit. So, let's suppose the Tigers go undefeated for the rest of the season and the Yankees fail to win another game.*

*The problem with this scenario is that New York still has 8 games left with Boston. If the Red Sox win all of these games, they will end the season with at least 77 wins putting them ahead of the Tigers. Thus, the only way for Detroit to even have a chance of finishing in first place, is for New York to win exactly one of the 8 games with Boston and lose all their other games. Meanwhile, the Sox must loss all the games they play agains teams other than New York. This puts them in a 3-way tie for first place....*

*Now let's look at what happens to the Orioles and Blue Jays in our scenario. Baltimore has 2 games left with Boston and 3 with New York. So, if everything happens as described above, the Orioles will finish with at least 76 wins. So, Detroit can catch Baltimore only if the Orioles lose all their games to teams other than New York and Boston. In particular, this means that Baltimore must lose all 7 of its remaining games with Toronto. The Blue Jays also have 7 games left with the Yankees and we have already seen that for Detroit to finish in first place, Toronto must will all of these games. But if that happens, the Blue Jays will win at least 14 more games giving them at final record of 77 and 85 or better which means they will finish ahead of the Tigers. So, no matter what happens from this point in the season on, Detroit can not finish in first place in the American League East.*

There has to be a better way to figure this out!

Here is a more abstract formulation of the problem. Our input consists of two arrays  $W[1..n]$  and  $G[1..n, 1..n]$ , where  $W[i]$  is the number of games team  $i$  has already won, and  $G[i, j]$  is the number of upcoming games between teams  $i$  and  $j$ . We want to determine whether team  $n$  can end the season with the most wins (possibly tied with other teams).<sup>3</sup>

We model this question as an assignment problem: We want to **assign** a winner to each game, so that team  $n$  comes in first place. We have an assignment problem! Let  $R[i] = \sum_j G[i, j]$  denote the number of remaining games for team  $i$ . We will assume that team  $n$  wins all  $R[n]$  of its remaining games. Then team  $n$  can come in first place if and only if every other team  $i$  wins at most  $W[n] + R[n] - W[i]$  of its  $R[i]$  remaining games.

Since we want to **assign** winning teams to games, we start by building a bipartite graph, whose nodes represent the games and the teams. We have  $\binom{n}{2}$  game nodes  $g_{i,j}$ , one for each pair  $1 \leq i < j \leq n$ , and  $n - 1$  team nodes  $t_i$ , one for each  $1 \leq i < n$ . For each pair  $i, j$ , we add edges  $g_{i,j} \rightarrow t_i$  and  $g_{i,j} \rightarrow t_j$  with infinite capacity. We add a source vertex  $s$  and edges  $s \rightarrow g_{i,j}$  with capacity  $G[i, j]$  for each pair  $i, j$ . Finally, we add a target node  $t$  and edges  $t_i \rightarrow t$  with capacity  $W[n] - W[i] + R[n]$  for each team  $i$ .

<sup>2</sup>Both the example and this argument are taken from <http://riot.ieor.berkeley.edu/~baseball/detroit.html>.

<sup>3</sup>We assume here that no games end in a tie (always true for Major League Baseball), and that every game is actually played (not always true).

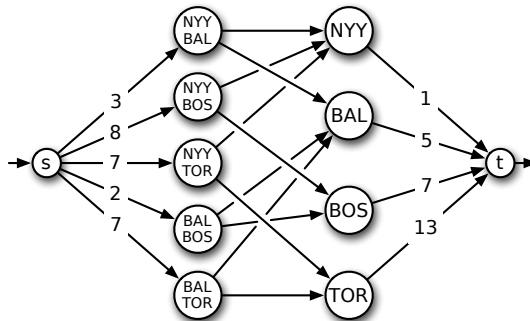
**Theorem:** Team  $n$  can end the season in first place if and only if there is a feasible flow in this graph that saturates every edge leaving  $s$ .

**Proof:** Suppose it is possible for team  $n$  to end the season in first place. Then every team  $i < n$  wins at most  $W[n] + R[n] - W[i]$  of the remaining games. For each game between team  $i$  and team  $j$  that team  $i$  wins, add one unit of flow along the path  $s \rightarrow g_{i,j} \rightarrow t_i \rightarrow t$ . Because there are exactly  $G[i, j]$  games between teams  $i$  and  $j$ , every edge leaving  $s$  is saturated. Because each team  $i$  wins at most  $W[n] + R[n] - W[i]$  games, the resulting flow is feasible.

Conversely, Let  $f$  be a feasible flow that saturates every edge out of  $s$ . Suppose team  $i$  wins exactly  $f(g_{i,j} \rightarrow t_i)$  games against team  $j$ , for all  $i$  and  $j$ . Then teams  $i$  and  $j$  play  $f(g_{i,j} \rightarrow t_i) + f(g_{i,j} \rightarrow t_j) = f(s \rightarrow g_{i,j}) = G[i, j]$  games, so every upcoming game is played. Moreover, each team  $i$  wins a total of  $\sum_j f(g_{i,j} \rightarrow t_i) = f(t_i \rightarrow t) \leq W[n] + R[n] - W[i]$  upcoming games, and therefore at most  $W[n] + R[n]$  games overall. Thus, if team  $n$  win all their upcoming games, they end the season in first place.  $\square$

So, to decide whether our favorite team can win, we construct the flow network, compute a maximum flow, and report whether than maximum flow saturates the edges leaving  $s$ . The flow network has  $O(n^2)$  vertices and  $O(n^2)$  edges, and it can be constructed in  $O(n^2)$  time. Using Dinitz's algorithm, we can compute the maximum flow in  $O(VE^2) = O(n^6)$  time.

The graph derived from the 1996 American League East standings is shown below. The total capacity of the edges leaving  $s$  is 27 (there are 27 remaining games), but the total capacity of the edges entering  $t$  is only 26. So the maximum flow has value at most 26, which means that Detroit is mathematically eliminated.

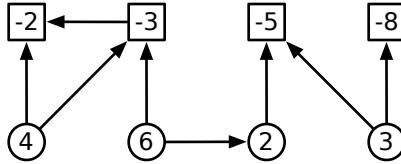


The flow graph for the 1996 American League East standings. Unlabeled edges have infinite capacity.

## 23.6 Project Selection

In our final example, suppose we are given a set of  $n$  projects that we could possibly perform; for simplicity, we identify each project by an integer between 1 and  $n$ . Some projects cannot be started until certain other projects are completed. This set of dependencies is described by a directed acyclic graph, where an edge  $i \rightarrow j$  indicates that project  $i$  depends on project  $j$ . Finally, each project  $i$  has an associated profit  $p_i$  which is given to us if the project is completed; however, some projects have negative profits, which we interpret as positive costs. We can choose to finish any subset  $X$  of the projects that includes all its dependents; that is, for every project  $x \in X$ , every project that  $x$  depends on is also in  $X$ . Our goal is to find a valid subset of the projects whose total profit is as large as possible. In particular, if all of the jobs have negative profit, the correct answer is to do nothing.

At a high level, our task to partition the projects into two subsets  $S$  and  $T$ , the jobs we *Select* and the jobs we *Turn down*. So intuitively, we'd like to model our problem as a minimum cut problem in a certain



A dependency graph for a set of projects. Circles represent profitable projects; squares represent costly projects.

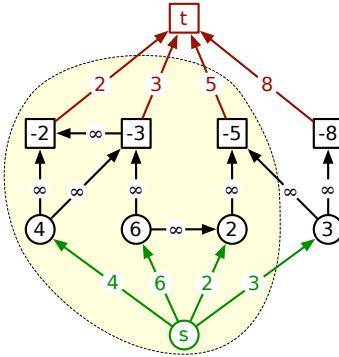
graph. But in which graph? How do we enforce prerequisites? We want to *maximize* profit, but we only know how to find *minimum* cuts. And how do we convert negative profits into positive capacities?

We define a new graph  $G$  by adding a source vertex  $s$  and a target vertex  $t$  to the dependency graph, with an edge  $s \rightarrow j$  for every profitable job (with  $p_j > 0$ ), and an edge  $i \rightarrow t$  for every costly job (with  $p_i < 0$ ). Intuitively, we can think of  $s$  as a new job (“To the bank!”) with profit/cost 0 that we must perform last. We assign edge capacities as follows:

- $c(s \rightarrow j) = p_j$  for every profitable job  $j$ ;
- $c(i \rightarrow t) = -p_i$  for every costly job  $i$ ;
- $c(i \rightarrow j) = \infty$  for every dependency edge  $i \rightarrow j$ .

All edge-capacities are positive, so this is a legal input to the maximum cut problem.

Now consider an  $(s, t)$ -cut  $(S, T)$  in  $G$ . If the capacity  $\|S, T\|$  is finite, then for every dependency edge  $i \rightarrow j$ , projects  $i$  and  $j$  are on the same side of the cut, which implies that  $S$  is a valid solution. Moreover, we claim that selecting the jobs in  $S$  earns us a total profit of  $C - \|S, T\|$ , where  $C$  is the sum of all the positive profits. This claim immediately implies that we can *maximize* our total profit by computing a *minimum* cut in  $G$ .



The flow network for the example dependency graph, along with its minimum cut.  
The cut has capacity 13 and  $C = 15$ , so the total profit for the selected jobs is 2.

We prove our key claim as follows. For any subset  $A$  of projects, we define three functions:

$$\begin{aligned} cost(A) &:= \sum_{i \in A: p_i < 0} -p_i = \sum_{i \in A} c(i \rightarrow t) \\ benefit(A) &:= \sum_{j \in A: p_j > 0} p_j = \sum_{j \in A} c(s \rightarrow j) \\ profit(A) &:= \sum_{i \in A} p_i = benefit(A) - cost(A). \end{aligned}$$

By definition,  $C = benefit(S) + benefit(T)$ . Because the cut  $(S, T)$  has finite capacity, only edges of the form  $s \rightarrow j$  and  $i \rightarrow t$  can cross the cut. By construction, every edge  $s \rightarrow j$  points to a profitable job and each

edge  $i \rightarrow t$  points from a costly job. Thus,  $\|S, T\| = cost(S) + benefit(T)$ . We immediately conclude that  $C - \|S, T\| = benefit(S) - cost(S) = profit(S)$ , as claimed.

## Exercises

1. Given an undirected graph  $G = (V, E)$ , with three vertices  $u$ ,  $v$ , and  $w$ , describe and analyze an algorithm to determine whether there is a path from  $u$  to  $w$  that passes through  $v$ .
2. Let  $G = (V, E)$  be a directed graph where for each vertex  $v$ , the in-degree and out-degree of  $v$  are equal. Let  $u$  and  $v$  be two vertices  $G$ , and suppose  $G$  contains  $k$  edge-disjoint paths from  $u$  to  $v$ . Under these conditions, must  $G$  also contain  $k$  edge-disjoint paths from  $v$  to  $u$ ? Give a proof or a counterexample with explanation.
3. A *cycle cover* of a given directed graph  $G = (V, E)$  is a set of vertex-disjoint cycles that cover all the vertices. Describe and analyze an efficient algorithm to find a cycle cover for a given graph, or correctly report that no cycle cover exists. [*Hint: Use bipartite matching!*]
4. Consider a directed graph  $G = (V, E)$  with multiple source vertices  $s_1, s_2, \dots, s_\sigma$  and multiple target vertices  $t_1, t_2, \dots, t_\tau$ , where no vertex is both a source and a target. A *multiterminal flow* is a function  $f : E \rightarrow \mathbb{R}_{\geq 0}$  that satisfies the flow conservation constraint at every vertex that is neither a source nor a target. The value  $|f|$  of a multiterminal flow is the total excess flow out of *all* the source vertices:

$$|f| := \sum_{i=1}^{\sigma} \left( \sum_w f(s_i \rightarrow w) - \sum_u f(u \rightarrow s_i) \right)$$

As usual, we are interested in finding flows with maximum value, subject to capacity constraints on the edges. (In particular, we don't care how much flow moves from any particular source to any particular target.)

- (a) Consider the following algorithm for computing multiterminal flows. The variables  $f$  and  $f'$  represent flow functions. The subroutine  $\text{MAXFLOW}(G, s, t)$  solves the standard maximum flow problem with source  $s$  and target  $t$ .

```
MAXMULTIFLOW( $G, s[1..\sigma], t[1..\tau]$ ):
   $f \leftarrow 0$                                 «Initialize the flow»
  for  $i \leftarrow 1$  to  $\sigma$ 
    for  $j \leftarrow 1$  to  $\tau$ 
       $f' \leftarrow \text{MAXFLOW}(G_f, s[i], t[j])$ 
       $f \leftarrow f + f'$                          «Update the flow»
  return  $f$ 
```

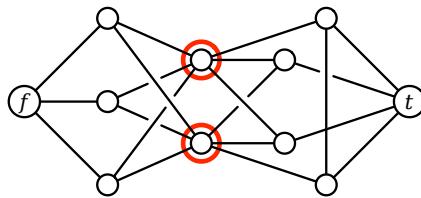
Prove that this algorithm correctly computes a maximum multiterminal flow in  $G$ .

- (b) Describe a more efficient algorithm to compute a maximum multiterminal flow in  $G$ .
5. The Island of Sodor is home to a large number of towns and villages, connected by an extensive rail network. Recently, several cases of a deadly contagious disease (either swine flu or zombies; reports are unclear) have been reported in the village of Ffarquhar. The controller of the Sodor railway plans to close down certain railway stations to prevent the disease from spreading to

Tidmouth, his home town. No trains can pass through a closed station. To minimize expense (and public notice), he wants to close down as few stations as possible. However, he cannot close the Ffarquhar station, because that would expose him to the disease, and he cannot close the Tidmouth station, because then he couldn't visit his favorite pub.

Describe and analyze an algorithm to find the minimum number of stations that must be closed to block all rail travel from Ffarquhar to Tidmouth. The Sodor rail network is represented by an undirected graph, with a vertex for each station and an edge for each rail connection between two stations. Two special vertices  $f$  and  $t$  represent the stations in Ffarquhar and Tidmouth.

For example, given the following input graph, your algorithm should return the number 2.



6. The UIUC Faculty Senate has decided to convene a committee to determine whether Chief Illiniwek should become the official *maseot symbol* of the University of Illinois Global Campus. Exactly one faculty member must be chosen from each academic department to serve on this committee. Some faculty members have appointments in multiple departments, but each committee member will represent only one department. For example, if Prof. Blagojevich is affiliated with both the Department of Corruption and the Department of Stupidity, and he is chosen as the Stupidity representative, then someone else must represent Corruption. Finally, University policy requires that any committee on virtual *maseots symbols* must contain the same number of assistant professors, associate professors, and full professors. Fortunately, the number of departments is a multiple of 3.

Describe an efficient algorithm to select the membership of the Global Illiniwek Committee. Your input is a list of all UIUC faculty members, their ranks (assistant, associate, or full), and their departmental affiliation(s). There are  $n$  faculty members and  $3k$  departments.

7. You're organizing the First Annual UIUC Computer Science 72-Hour Dance Exchange, to be held all day Friday, Saturday, and Sunday. Several 30-minute sets of music will be played during the event, and a large number of DJs have applied to perform. You need to hire DJs according to the following constraints.

- Exactly  $k$  sets of music must be played each day, and thus  $3k$  sets altogether.
- Each set must be played by a single DJ in a consistent music genre (ambient, bubblegum, dubstep, horrorcore, hyphy, trip-hop, Nitzhonot, Kwaito, J-pop, Nashville country, ...).
- Each genre must be played at most once per day.
- Each candidate DJ has given you a list of genres they are willing to play.
- Each DJ can play at most three sets during the entire event.

Suppose there are  $n$  candidate DJs and  $g$  different musical genres available. Describe and analyze an efficient algorithm that either assigns a DJ and a genre to each of the  $3k$  sets, or correctly reports that no such assignment is possible.

8. The University of Southern North Dakota at Hoople has hired you to write an algorithm to schedule their final exams. Each semester, USNDH offers  $n$  different classes. There are  $r$  different rooms on campus and  $t$  different time slots in which exams can be offered. You are given two arrays  $E[1..n]$  and  $S[1..r]$ , where  $E[i]$  is the number of students enrolled in the  $i$ th class, and  $S[j]$  is the number of seats in the  $j$ th room. At most one final exam can be held in each room during each time slot. Class  $i$  can hold its final exam in room  $j$  only if  $E[i] < S[j]$ .

Describe and analyze an efficient algorithm to assign a room and a time slot to each class (or report correctly that no such assignment is possible).

9. Suppose you are running a web site that is visited by the same set of people every day. Each visitor claims membership in one or more *demographic groups*; for example, a visitor might describe himself as male, 35–45 years old, a father, a resident of Illinois, an academic, a blogger, and a fan of Joss Whedon.<sup>4</sup> Your site is supported by advertisers. Each advertiser has told you which demographic groups should see its ads and how many of its ads you must show each day. Altogether, there are  $n$  visitors,  $k$  demographic groups, and  $m$  advertisers.

Describe an efficient algorithm to determine, given all the data described in the previous paragraph, whether you can show each visitor exactly *one* ad per day, so that every advertiser has its desired number of ads displayed, and every ad is seen by someone in an appropriate demographic group.

10. Suppose we are given an array  $A[1..m][1..n]$  of non-negative real numbers. We want to *round*  $A$  to an integer matrix, by replacing each entry  $x$  in  $A$  with either  $\lfloor x \rfloor$  or  $\lceil x \rceil$ , without changing the sum of entries in any row or column of  $A$ . For example:

$$\begin{bmatrix} 1.2 & 3.4 & 2.4 \\ 3.9 & 4.0 & 2.1 \\ 7.9 & 1.6 & 0.5 \end{bmatrix} \longrightarrow \begin{bmatrix} 1 & 4 & 2 \\ 4 & 4 & 2 \\ 8 & 1 & 1 \end{bmatrix}$$

- (a) Describe an efficient algorithm that either rounds  $A$  in this fashion, or reports correctly that no such rounding is possible.
- \*(b) Suppose we are guaranteed that none of the entries in the input matrix  $A$  are integers. Describe an even faster algorithm that either rounds  $A$  or reports correctly that no such rounding is possible. For full credit, your algorithm must run in  $O(mn)$  time. [Hint: **Don't use flows.**]

11. *Ad-hoc networks* are made up of low-powered wireless devices. In principle<sup>5</sup>, these networks can be used on battlefields, in regions that have recently suffered from natural disasters, and in other hard-to-reach areas. The idea is that a large collection of cheap, simple devices could be distributed through the area of interest (for example, by dropping them from an airplane); the devices would then automatically configure themselves into a functioning wireless network.

These devices can communicate only within a limited range. We assume all the devices are identical; there is a distance  $D$  such that two devices can communicate if and only if the distance between them is at most  $D$ .

---

<sup>4</sup>Har har har! Mine is an evil laugh! Now die!

<sup>5</sup>but not really in practice

We would like our ad-hoc network to be reliable, but because the devices are cheap and low-powered, they frequently fail. If a device detects that it is likely to fail, it should transmit its information to some other *backup* device within its communication range. We require each device  $x$  to have  $k$  potential backup devices, all within distance  $D$  of  $x$ ; we call these  $k$  devices the *backup set* of  $x$ . Also, we do not want any device to be in the backup set of too many other devices; otherwise, a single failure might affect a large fraction of the network.

So suppose we are given the communication radius  $D$ , parameters  $b$  and  $k$ , and an array  $d[1..n, 1..n]$  of distances, where  $d[i, j]$  is the distance between device  $i$  and device  $j$ . Describe an algorithm that either computes a backup set of size  $k$  for each of the  $n$  devices, such that no device appears in more than  $b$  backup sets, or reports (correctly) that no good collection of backup sets exists.

- \*12. A *rooted tree* is a directed acyclic graph, in which every vertex has exactly one incoming edge, except for the *root*, which has no incoming edges. Equivalently, a rooted tree consists of a root vertex, which has edges pointing to the roots of zero or more smaller rooted trees. Describe a polynomial-time algorithm to compute, given two rooted trees  $A$  and  $B$ , the largest common rooted subtree of  $A$  and  $B$ .

[Hint: Let  $LCS(u, v)$  denote the largest common subtree whose root in  $A$  is  $u$  and whose root in  $B$  is  $v$ . Your algorithm should compute  $LCS(u, v)$  for all vertices  $u$  and  $v$  using dynamic programming. This would be easy if every vertex had  $O(1)$  children, and still straightforward if the children of each node were ordered from left to right and the common subtree had to respect that ordering. But for unordered trees with large degree, you need another trick to combine recursive subproblems efficiently. Don't waste your time trying to reduce the polynomial running time.]

"Who are you?" said Lunkwill, rising angrily from his seat. "What do you want?"  
 "I am Majikthise!" announced the older one.  
 "And I demand that I am Vroomfondel!" shouted the younger one.  
 Majikthise turned on Vroomfondel. "It's alright," he explained angrily, "you don't need to demand that."  
 "Alright!" bawled Vroomfondel banging on a nearby desk. "I am Vroomfondel, and that is not a demand, that is a solid fact! What we demand is solid facts!"  
 "No we don't!" exclaimed Majikthise in irritation. "That is precisely what we don't demand!"  
 Scarcely pausing for breath, Vroomfondel shouted, "We don't demand solid facts! What we demand is a total absence of solid facts. I demand that I may or may not be Vroomfondel!"

— Douglas Adams, *The Hitchhiker's Guide to the Galaxy* (1979)

## \*18 Extensions of Maximum Flow

### 18.1 Maximum Flows with Edge Demands

Now suppose each directed edge  $e$  in has both a capacity  $c(e)$  and a *demand*  $d(e) \leq c(e)$ , and we want a flow  $f$  of maximum value that satisfies  $d(e) \leq f(e) \leq c(e)$  at every edge  $e$ . We call a flow that satisfies these constraints a *feasible* flow. In our original setting, where  $d(e) = 0$  for every edge  $e$ , the zero flow is feasible; however, in this more general setting, even determining whether a feasible flow exists is a nontrivial task.

Perhaps the easiest way to find a feasible flow (or determine that none exists) is to reduce the problem to a standard maximum flow problem, as follows. The input consists of a directed graph  $G = (V, E)$ , nodes  $s$  and  $t$ , demand function  $d: E \rightarrow \mathbb{R}$ , and capacity function  $c: E \rightarrow \mathbb{R}$ . Let  $D$  denote the sum of all edge demands in  $G$ :

$$D := \sum_{u \rightarrow v \in E} d(u \rightarrow v).$$

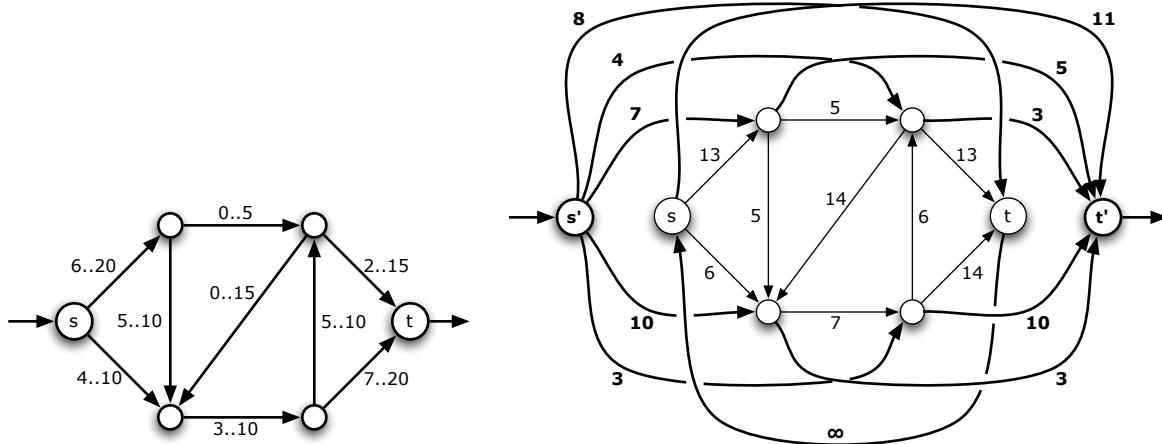
We construct a new graph  $G' = (V', E')$  from  $G$  by adding new source and target vertices  $s'$  and  $t'$ , adding edges from  $s'$  to each vertex in  $V$ , adding edges from each vertex in  $V$  to  $t'$ , and finally adding an edge from  $t$  to  $s$ . We also define a new capacity function  $c': E' \rightarrow \mathbb{R}$  as follows:

- For each vertex  $v \in V$ , we set  $c'(s' \rightarrow v) = \sum_{u \in V} d(u \rightarrow v)$  and  $c'(v \rightarrow t') = \sum_{w \in V} d(v \rightarrow w)$ .
- For each edge  $u \rightarrow v \in E$ , we set  $c'(u \rightarrow v) = c(u \rightarrow v) - d(u \rightarrow v)$ .
- Finally, we set  $c'(t \rightarrow s) = \infty$ .

Intuitively, we construct  $G'$  by replacing any edge  $u \rightarrow v$  in  $G$  with three edges: an edge  $u \rightarrow v$  with capacity  $c(u \rightarrow v) - d(u \rightarrow v)$ , an edge  $s' \rightarrow v$  with capacity  $d(u \rightarrow v)$ , and an edge  $u \rightarrow t'$  with capacity  $d(u \rightarrow v)$ . If this construction produces multiple edges from  $s'$  to the same vertex  $v$  (or to  $t'$  from the same vertex  $v$ ), we merge them into a single edge with the same total capacity.

In  $G'$ , the total capacity out of  $s'$  and the total capacity into  $t'$  are both equal to  $D$ . We call a flow with value exactly  $D$  a *saturating* flow, since it saturates all the edges leaving  $s'$  or entering  $t'$ . If  $G'$  has a saturating flow, it must be a maximum flow, so we can find it using any max-flow algorithm.

**Lemma 1.**  $G$  has a feasible  $(s, t)$ -flow if and only if  $G'$  has a saturating  $(s', t')$ -flow.



A flow network  $G$  with demands and capacities (written  $d \dots c$ ), and the transformed network  $G'$ .

**Proof:** Let  $f : E \rightarrow \mathbb{R}$  be a feasible  $(s, t)$ -flow in the original graph  $G$ . Consider the following function  $f' : E' \rightarrow \mathbb{R}$ :

$$\begin{aligned} f'(u \rightarrow v) &= f(u \rightarrow v) - d(u \rightarrow v) && \text{for all } u \rightarrow v \in E \\ f'(s' \rightarrow v) &= \sum_{u \in V} d(u \rightarrow v) && \text{for all } v \in V \\ f'(v \rightarrow t') &= \sum_{w \in V} d(v \rightarrow w) && \text{for all } v \in V \\ f'(t \rightarrow s) &= |f| \end{aligned}$$

We easily verify that  $f'$  is a saturating  $(s', t')$ -flow in  $G$ . The admissibility of  $f$  implies that  $f(e) \geq d(e)$  for every edge  $e \in E$ , so  $f'(e) \geq 0$  everywhere. Admissibility also implies  $f(e) \leq c(e)$  for every edge  $e \in E$ , so  $f'(e) \leq c'(e)$  everywhere. Tedious algebra implies that

$$\sum_{u \in V'} f'(u \rightarrow v) = \sum_{w \in V'} f(v \rightarrow w)$$

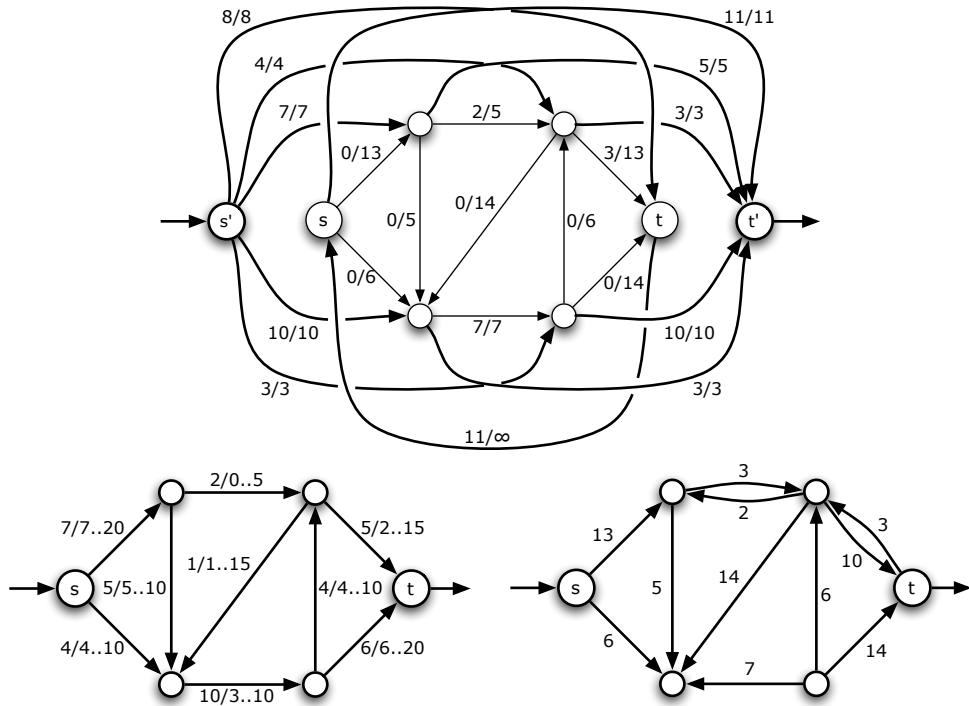
for every vertex  $v \in V$  (including  $s$  and  $t$ ). Thus,  $f'$  is a legal  $(s', t')$ -flow, and every edge out of  $s'$  or into  $t'$  is clearly saturated. Intuitively,  $f'$  diverts  $d(u \rightarrow v)$  units of flow from  $u$  directly to the new target  $t'$ , and injects the same amount of flow into  $v$  directly from the new source  $s'$ .

The same tedious algebra implies that for any saturating  $(s', t')$ -flow  $f' : E' \rightarrow \mathbb{R}$  for  $G'$ , the function  $f = f'|_E + d$  is a feasible  $(s, t)$ -flow in  $G$ .  $\square$

Thus, we can compute a feasible  $(s, t)$ -flow for  $G$ , if one exists, by searching for a maximum  $(s', t')$ -flow in  $G'$  and checking that it is saturating. Once we've found a feasible  $(s, t)$ -flow in  $G$ , we can transform it into a maximum flow using an augmenting-path algorithm, but with one small change. To ensure that *every* flow we consider is feasible, we must redefine the residual capacity of an edge as follows:

$$c_f(u \rightarrow v) = \begin{cases} c(u \rightarrow v) - f(u \rightarrow v) & \text{if } u \rightarrow v \in E, \\ f(v \rightarrow u) - d(v \rightarrow u) & \text{if } v \rightarrow u \in E, \\ 0 & \text{otherwise.} \end{cases}$$

Otherwise, the algorithm is unchanged. If we use the Dinitz/Edmonds-Karp fat-pipe algorithm, we get an overall running time of  $O(VE^2)$ .



A saturating flow  $f'$  in  $G'$ , the corresponding feasible flow  $f$  in  $G$ , and the corresponding residual network  $G_f$ .

## 18.2 Node Supplies and Demands

Another useful variant to consider allows flow to be injected or extracted from the flow network at vertices other than  $s$  or  $t$ . Let  $x: (V \setminus \{s, t\}) \rightarrow \mathbb{R}$  be an excess function describing how much flow is to be injected (or extracted if the value is negative) at each vertex. We now want a maximum ‘flow’ that satisfies the variant balance condition

$$\sum_{u \in V} f(u \rightarrow v) - \sum_{w \in V} f(v \rightarrow w) = x(v)$$

for every node  $v$  except  $s$  and  $t$ , or prove that no such flow exists. As above, call such a function  $f$  a *feasible flow*.

As for flows with edge demands, the only real difficulty in finding a maximum flow under these modified constraints is finding a feasible flow (if one exists). We can reduce this problem to a standard max-flow problem, just as we did for edge demands.

To simplify the transformation, let us assume without loss of generality that the total excess in the network is zero:  $\sum_v x(v) = 0$ . If the total excess is positive, we add an infinite capacity edge  $t \rightarrow \tilde{t}$ , where  $\tilde{t}$  is a new target node, and set  $x(\tilde{t}) = -\sum_v x(v)$ . Similarly, if the total excess is negative, we add an infinite capacity edge  $\tilde{s} \rightarrow s$ , where  $\tilde{s}$  is a new source node, and set  $x(s) = -\sum_v x(v)$ . In both cases, every feasible flow in the modified graph corresponds to a feasible flow in the original graph.

As before, we modify  $G$  to obtain a new graph  $G'$  by adding a new source  $s'$ , a new target  $t'$ , an infinite-capacity edge  $t \rightarrow s$  from the old target to the old source, and several edges from  $s'$  and to  $t'$ . Specifically, for each vertex  $v$ , if  $x(v) > 0$ , we add a new edge  $s' \rightarrow v$  with capacity  $x(v)$ , and if  $x(v) < 0$ , we add an edge  $v \rightarrow t'$  with capacity  $-x(v)$ . As before, we call an  $(s', t')$ -flow in  $G'$  *saturating* if every edge leaving  $s'$  or entering  $t'$  is saturated; any saturating flow is a maximum flow. It is easy to check that saturating flows in  $G'$  are in direct correspondence with feasible flows in  $G$ ; we leave details as an exercise (hint, hint).

Similar reductions allow us to several other variants of the maximum flow problem using the same path-augmentation techniques. For example, we could associate capacities and demands with the vertices instead of (or in addition to) the edges, as well as a *range* of excesses with every vertex, instead of a single excess value.

### 18.3 Minimum-Cost Flows

Now imagine that each edge  $e$  in the network has both a capacity  $c(e)$  and a cost  $\$(e)$ . The cost function describes the cost of sending a unit of flow through the edges; thus, the cost any flow  $f$  is defined as follows:

$$\$(f) = \sum_{e \in E} \$(e) \cdot f(e).$$

The *minimum-cost maximum-flow* problem is to compute a maximum flow of minimum cost. If the network has only one maximum flow, that's what we want, but if there is more than one maximum flow, we want the maximum flow whose cost is as small as possible. Costs can either be positive, negative, or zero. However, if an edge  $u \rightarrow v$  and its reversal  $v \rightarrow u$  both appear in the graph, their costs must sum to zero:  $\$(u \rightarrow v) = -\$(v \rightarrow u)$ . Otherwise, we could make an infinite profit by pushing flow back and forth along the edge!

Each augmentation step in the standard Ford-Fulkerson algorithm both increases the value of the flow and changes its cost. If the total cost of the augmenting path is positive, the cost of the flow decreases; conversely, if the total cost of the augmenting path is negative, the cost of the flow increases. We can also change the cost of the flow without changing its value, by augmenting along a directed *cycle* in the residual graph. Again, augmenting along a negative-cost cycle decreases the cost of the flow, and augmenting along a positive-cost cycle increases the cost of the flow.

It follows immediately that a flow  $f$  is a minimum-cost maximum flow in  $G$  if and only if the residual graph  $G_f$  has no directed paths from  $s$  to  $t$  and no negative-cost cycles.

We can compute a min-cost max-flow using the so-called *cycle cancelling* algorithm first proposed by Morton Klein in 1967. The algorithm has two phases; in the first, we compute an arbitrary maximum flow  $f$ , using any method we like. The second phase repeatedly decreases the cost of  $f$ , by augmenting  $f$  along a negative-cost cycle in the residual graph  $G_f$ , until no such cycle exists. As in Ford-Fulkerson, the amount of flow we push around each cycle is equal to the minimum residual capacity of any edge on the cycle.

In each iteration of the second phase, we can use a modification of Shimbels shortest path algorithm (often called “Bellman-Ford”) to find a negative-cost cycle in  $O(VE)$  time. To bound the number of iterations in the second phase, we assume that both the capacity and the cost of each edge is an integer, and we define

$$C = \max_{e \in E} c(e) \quad \text{and} \quad D = \max_{e \in E} |\$(e)|.$$

The cost of any feasible flow is clearly between  $-ECD$  and  $ECD$ , and each augmentation step decreases the cost of the flow by a positive integer, and therefore by at least 1. We conclude that the second phase requires at most  $2ECD$  iterations, and therefore runs in  $O(VE^2CD)$  time. As with the raw Ford-Fulkerson algorithm, this running time is exponential in the complexity of the input, and it may never terminate if the capacities and/or costs are irrational.

Like Ford-Fulkerson, more careful choices of which cycle to cancel can lead to more efficient algorithms. Unfortunately, some obvious choices are NP-hard to compute, including the cycle with most negative cost and the negative cycle with the fewest edges. In the late 1980s, Andrew Goldberg and Bob Tarjan developed a min-cost flow algorithm that repeatedly cancels the so-called *minimum-mean cycle*, which is the cycle whose *average cost per edge* is smallest. By combining an algorithm of Karp to compute

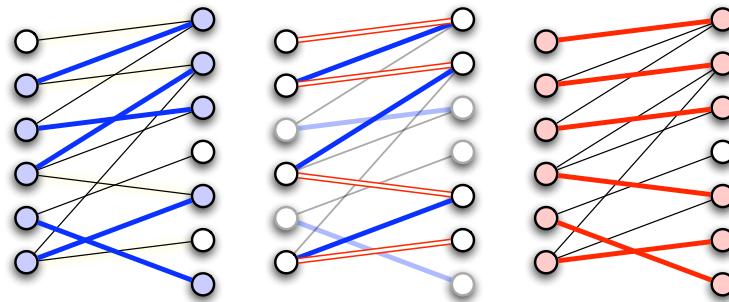
minimum-mean cycles in  $O(EV)$  time, efficient dynamic tree data structures, and other sophisticated techniques that are (unfortunately) beyond the scope of this class, their algorithm achieves a running time of  $O(E^2V \log^2 V)$ . The fastest min-cost max-flow algorithm currently known,<sup>1</sup> due to James Orlin, reduces the problem to  $O(E \log V)$  iterations of Dijkstra's shortest-path algorithm; Orlin's algorithm runs in  $O(E^2 \log V + EV \log^2 V)$  time.

## 18.4 Maximum-Weight Matchings

Recall from the previous lecture that we can find a maximum-cardinality matching in any bipartite graph in  $O(VE)$  time by reduction to the standard maximum flow problem.

Now suppose the input graph has *weighted* edges, and we want to find the matching with maximum total weight. Given a bipartite graph  $G = (U \times W, E)$  and a non-negative weight function  $w: E \rightarrow \mathbb{R}$ , the goal is to compute a matching  $M$  whose total weight  $w(M) = \sum_{uw \in M} w(uw)$  is as large as possible. Max-weight matchings can't be found directly using standard max-flow algorithms<sup>2</sup>, but we can modify the algorithm for maximum-cardinality matchings described above.

It will be helpful to reinterpret the behavior of our earlier algorithm directly in terms of the original bipartite graph instead of the derived flow network. Our algorithm maintains a matching  $M$ , which is initially empty. We say that a vertex is *matched* if it is an endpoint of an edge in  $M$ . At each iteration, we find an *alternating path*  $\pi$  that starts and ends at unmatched vertices and alternates between edges in  $E \setminus M$  and edges in  $M$ . Equivalently, let  $G_M$  be the directed graph obtained by orienting every edge in  $M$  from  $W$  to  $U$ , and every edge in  $E \setminus M$  from  $U$  to  $W$ . An alternating path is just a directed path in  $G_M$  between two unmatched vertices. Any alternating path has odd length and has exactly one more edge in  $E \setminus M$  than in  $M$ . The iteration ends by setting  $M \leftarrow M \oplus \pi$ , thereby increasing the number of edges in  $M$  by one. The max-flow/min-cut theorem implies that when there are no more alternating paths,  $M$  is a maximum matching.



A matching  $M$  with 5 edges, an alternating path  $\pi$ , and the augmented matching  $M \oplus \pi$  with 6 edges.

If the edges of  $G$  are weighted, we only need to make two changes to the algorithm. First, instead of looking for an arbitrary alternating path at each iteration, we look for the alternating path  $\pi$  such that  $M \oplus \pi$  has largest weight. Suppose we weight the edges in the residual graph  $G_M$  as follows:

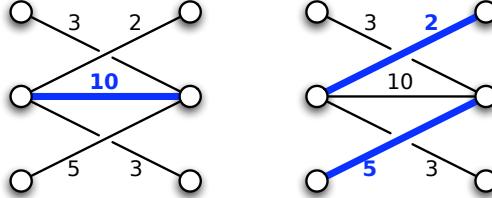
$$\begin{aligned} w'(u \rightarrow w) &= -w(uw) && \text{for all } uw \notin M \\ w'(w \rightarrow u) &= w(uw) && \text{for all } uw \in M \end{aligned}$$

We now have  $w(M \oplus \pi) = w(M) - w'(\pi)$ . **Thus, the correct augmenting path  $\pi$  must be the directed path in  $G_M$  with minimum total residual weight  $w'(\pi)$ .** Second, because the matching with the

<sup>1</sup>at least, among algorithms whose running times do not depend on  $C$  and  $D$

<sup>2</sup>However, max-flow algorithms can be modified to compute maximum weighted flows, where every edge has both a capacity and a weight, and the goal is to maximize  $\sum_{u \rightarrow v} w(u \rightarrow v)f(u \rightarrow v)$ .

maximum *weight* may not be the matching with the maximum *cardinality*, we return the heaviest matching considered in *any* iteration of the algorithm.



A maximum-weight matching is not necessarily a maximum-cardinality matching.

Before we determine the running time of the algorithm, we need to check that it actually finds the maximum-weight matching. After all, it's a greedy algorithm, and greedy algorithms don't work unless you prove them into submission! Let  $M_i$  denote the maximum-weight matching in  $G$  with exactly  $i$  edges. In particular,  $M_0 = \emptyset$ , and the global maximum-weight matching is equal to  $M_i$  for some  $i$ . (The figure above show  $M_1$  and  $M_2$  for the same graph.) Let  $G_i$  denote the directed residual graph for  $M_i$ , let  $w_i$  denote the residual weight function for  $M_i$  as defined above, and let  $\pi_i$  denote the directed path in  $G_i$  such that  $w_i(\pi_i)$  is minimized. To simplify the proof, I will assume that there is a unique maximum-weight matching  $M_i$  of any particular size; this assumption can be enforced by applying a consistent tie-breaking rule. With this assumption in place, the correctness of our algorithm follows inductively from the following lemma.

**Lemma 2.** *If  $G$  contains a matching with  $i + 1$  edges, then  $M_{i+1} = M_i \oplus \pi_i$ .*

**Proof:** I will prove the equivalent statement  $M_{i+1} \oplus M_i = \pi_{i-1}$ . To simplify notation, call an edge in  $M_{i+1} \oplus M_i$  **red** if it is an edge in  $M_{i+1}$ , and **blue** if it is an edge in  $M_i$ .

The graph  $M_{i+1} \oplus M_i$  has maximum degree 2, and therefore consists of pairwise disjoint paths and cycles, each of which alternates between red and blue edges. Since  $G$  is bipartite, every cycle must have even length. The number of edges in  $M_{i+1} \oplus M_i$  is odd; specifically,  $M_{i+1} \oplus M_i$  has  $2i + 1 - 2k$  edges, where  $k$  is the number of edges that are in both matchings. Thus,  $M_{i+1} \oplus M_i$  contains an odd number of paths of odd length, some number of paths of even length, and some number of cycles of even length.

Let  $\gamma$  be a cycle in  $M_{i+1} \oplus M_i$ . Because  $\gamma$  has an equal number of edges from each matching,  $M_i \oplus \gamma$  is another matching with  $i$  edges. The total weight of this matching is exactly  $w(M_i) - w_i(\gamma)$ , which must be less than  $w(M_i)$ , so  $w_i(\gamma)$  must be positive. On the other hand,  $M_{i+1} \oplus \gamma$  is a matching with  $i + 1$  edges whose total weight is  $w(M_{i+1}) + w_i(\gamma) < w(M_{i+1})$ , so  $w_i(\gamma)$  must be negative! We conclude that no such cycle  $\gamma$  exists;  $M_{i+1} \oplus M_i$  consists entirely of disjoint paths.

Exactly the same reasoning implies that no path in  $M_{i+1} \oplus M_i$  has an even number of edges.

Finally, since the number of red edges in  $M_{i+1} \oplus M_i$  is one more than the number of blue edges, the number of paths that start with a red edge is exactly one more than the number of paths that start with a blue edge. The same reasoning as above implies that  $M_{i+1} \oplus M_i$  does not contain a blue-first path, because we can pair it up with a red-first path.

We conclude that  $M_{i+1} \oplus M_i$  consists of a single alternating path  $\pi$  whose first edge is red. Since  $w(M_{i+1}) = w(M_i) - w_i(\pi)$ , the path  $\pi$  must be the one with minimum weight  $w_i(\pi)$ .  $\square$

We can find the alternating path  $\pi_i$  using a single-source shortest path algorithm. Modify the residual graph  $G_i$  by adding zero-weight edges from a new source vertex  $s$  to every *unmatched* node in  $U$ , and from every *unmatched* node in  $W$  to a new target vertex  $t$ , exactly as in our unweighted matching algorithm. Then  $\pi_i$  is the shortest path from  $s$  to  $t$  in this modified graph. Since  $M_i$  is the maximum-weight matching with  $i$  vertices,  $G_i$  has no negative cycles, so this shortest path is well-defined.

We can compute the shortest path in  $G_i$  in  $O(VE)$  time using Shimbels algorithm, so the overall running time our algorithm is  $O(V^2E)$ .

The residual graph  $G_i$  has negative-weight edges, so we can't speed up the algorithm by replacing Shimbels algorithm with Dijkstra's. However, we can use a variant of Johnson's all-pairs shortest path algorithm to improve the running time to  $O(VE + V^2 \log V)$ . Let  $d_i(v)$  denote the distance from  $s$  to  $v$  in the residual graph  $G_i$ , using the distance function  $w_i$ . Let  $\tilde{w}_i$  denote the modified distance function  $\tilde{w}_i(u \rightarrow v) = d_{i-1}(u) + w_i(u \rightarrow v) - d_{i-1}(v)$ . As we argued in the discussion of Johnson's algorithm, shortest paths with respect to  $w_i$  are still shortest paths with respect to  $\tilde{w}_i$ . Moreover,  $\tilde{w}_i(u \rightarrow v) > 0$  for every edge  $u \rightarrow v$  in  $G_i$ :

- If  $u \rightarrow v$  is an edge in  $G_{i-1}$ , then  $w_i(u \rightarrow v) = w_{i-1}(u \rightarrow v)$  and  $d_{i-1}(v) \leq d_{i-1}(u) + w_{i-1}(u \rightarrow v)$ .
- If  $u \rightarrow v$  is not in  $G_{i-1}$ , then  $w_i(u \rightarrow v) = -w_{i-1}(v \rightarrow u)$  and  $v \rightarrow u$  is an edge in the shortest path  $\pi_{i-1}$ , so  $d_{i-1}(u) = d_{i-1}(v) + w_{i-1}(v \rightarrow u)$ .

Let  $\tilde{d}_i(v)$  denote the shortest path distance from  $s$  to  $v$  with respect to the distance function  $\tilde{w}_i$ . Because  $\tilde{w}_i$  is positive everywhere, we can quickly compute  $\tilde{d}_i(v)$  for all  $v$  using Dijkstra's algorithm. This gives us both the shortest alternating path  $\pi_i$  and the distances  $d_i(v) = \tilde{d}_i(v) + d_{i-1}(v)$  needed for the next iteration.

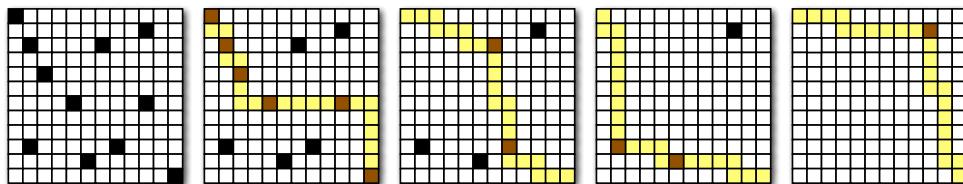
## Exercises

1. Suppose we are given a directed graph  $G = (V, E)$ , two vertices  $s$  and  $t$ , and a capacity function  $c: V \rightarrow \mathbb{R}^+$ . A flow  $f$  is *feasible* if the total flow into every vertex  $v$  is at most  $c(v)$ :

$$\sum_u f(u \rightarrow v) \leq c(v) \quad \text{for every vertex } v.$$

Describe and analyze an efficient algorithm to compute a feasible flow of maximum value.

2. Suppose we are given an  $n \times n$  grid, some of whose cells are marked; the grid is represented by an array  $M[1..n, 1..n]$  of booleans, where  $M[i, j] = \text{TRUE}$  if and only if cell  $(i, j)$  is marked. A *monotone* path through the grid starts at the top-left cell, moves only right or down at each step, and ends at the bottom-right cell. Our goal is to cover the marked cells with as few monotone paths as possible.



Greedyly covering the marked cells in a grid with four monotone paths.

- (a) Describe an algorithm to find a monotone path that covers the largest number of marked cells.
- (b) There is a natural greedy heuristic to find a small cover by monotone paths: If there are any marked cells, find a monotone path  $\pi$  that covers the largest number of marked cells, unmark any cells covered by  $\pi$  those marked cells, and recurse. Show that this algorithm does *not* always compute an optimal solution.

- (c) Describe and analyze an efficient algorithm to compute the smallest set of monotone paths that covers every marked cell.
3. Suppose we are given a set of boxes, each specified by their height, width, and depth in centimeters. All three side lengths of every box lie strictly between 10cm and 20cm. As you should expect, one box can be placed inside another if the smaller box can be rotated so that its height, width, and depth are respectively smaller than the height, width, and depth of the larger box. Boxes can be nested recursively. Call a box is *visible* if it is not inside another box.  
Describe and analyze an algorithm to nest the boxes so that the number of visible boxes is as small as possible.
4. Let  $G$  be a directed flow network whose edges have costs, but which contains no negative-cost cycles. Prove that one can compute a minimum-cost maximum flow in  $G$  using a variant of Ford-Fulkerson that repeatedly augments the  $(s, t)$ -path of *minimum total cost* in the current residual graph. What is the running time of this algorithm?

*The greatest flood has the soonest ebb;  
the sorest tempest the most sudden calm;  
the hottest love the coldest end; and  
from the deepest desire oftentimes ensues the deadliest hate.*

— Socrates

*Th' extremes of glory and of shame,  
Like east and west, become the same.*

— Samuel Butler, *Hudibras* Part II, Canto I (c. 1670)

*Extremes meet, and there is no better example  
than the haughtiness of humility.*

— Ralph Waldo Emerson, “Greatness”,  
in *Letters and Social Aims* (1876)

## \*25 Linear Programming

The maximum flow/minimum cut problem is a special case of a very general class of problems called *linear programming*. Many other optimization problems fall into this class, including minimum spanning trees and shortest paths, as well as several common problems in scheduling, logistics, and economics. Linear programming was used implicitly by Fourier in the early 1800s, but it was first formalized and applied to problems in economics in the 1930s by Leonid Kantorovich. Kantorovich’s work was hidden behind the Iron Curtain (where it was largely ignored) and therefore unknown in the West. Linear programming was rediscovered and applied to shipping problems in the early 1940s by Tjalling Koopmans. The first complete algorithm to solve linear programming problems, called the *simplex method*, was published by George Dantzig in 1947. Koopmans first proposed the name “linear programming” in a discussion with Dantzig in 1948. Kantorovich and Koopmans shared the 1975 Nobel Prize in Economics “for their contributions to the theory of optimum allocation of resources”. Dantzig did not; his work was apparently too pure. Koopmans wrote to Kantorovich suggesting that they refuse the prize in protest of Dantzig’s exclusion, but Kantorovich saw the prize as a vindication of his use of mathematics in economics, which his Soviet colleagues had written off as “a means for apologists of capitalism”.

A linear programming problem asks for a vector  $x \in \mathbb{R}^d$  that maximizes (or equivalently, minimizes) a given linear function, among all vectors  $x$  that satisfy a given set of linear inequalities. The general form of a linear programming problem is the following:

$$\begin{aligned} & \text{maximize } \sum_{j=1}^d c_j x_j \\ & \text{subject to } \sum_{j=1}^d a_{ij} x_j \leq b_i \quad \text{for each } i = 1 \dots p \\ & \quad \sum_{j=1}^d a_{ij} x_j = b_i \quad \text{for each } i = p+1 \dots p+q \\ & \quad \sum_{j=1}^d a_{ij} x_j \geq b_i \quad \text{for each } i = p+q+1 \dots n \end{aligned}$$

Here, the input consists of a matrix  $A = (a_{ij}) \in \mathbb{R}^{n \times d}$ , a column vector  $b \in \mathbb{R}^n$ , and a row vector  $c \in \mathbb{R}^d$ . Each coordinate of the vector  $x$  is called a *variable*. Each of the linear inequalities is called a *constraint*.

The function  $x \mapsto c \cdot x$  is called the *objective function*. I will always use  $d$  to denote the number of variables, also known as the *dimension* of the problem. The number of constraints is usually denoted  $n$ .

A linear programming problem is said to be in *canonical form*<sup>1</sup> if it has the following structure:

$$\begin{aligned} & \text{maximize} \sum_{j=1}^d c_j x_j \\ & \text{subject to } \sum_{j=1}^d a_{ij} x_j \leq b_i \quad \text{for each } i = 1..n \\ & \quad x_j \geq 0 \quad \text{for each } j = 1..d \end{aligned}$$

We can express this canonical form more compactly as follows. For two vectors  $x = (x_1, x_2, \dots, x_d)$  and  $y = (y_1, y_2, \dots, y_d)$ , the expression  $x \geq y$  means that  $x_i \geq y_i$  for every index  $i$ .

$$\begin{array}{l} \max c \cdot x \\ \text{s.t. } Ax \leq b \\ \quad x \geq 0 \end{array}$$

Any linear programming problem can be converted into canonical form as follows:

- For each variable  $x_j$ , add the equality constraint  $x_j = x_j^+ - x_j^-$  and the inequalities  $x_j^+ \geq 0$  and  $x_j^- \geq 0$ .
- Replace any equality constraint  $\sum_j a_{ij} x_j = b_i$  with two inequality constraints  $\sum_j a_{ij} x_j \geq b_i$  and  $\sum_j a_{ij} x_j \leq b_i$ .
- Replace any upper bound  $\sum_j a_{ij} x_j \geq b_i$  with the equivalent lower bound  $\sum_j -a_{ij} x_j \leq -b_i$ .

This conversion potentially double the number of variables and the number of constraints; fortunately, it is rarely necessary in practice.

Another useful format for linear programming problems is *slack form*<sup>2</sup>, in which every inequality is of the form  $x_j \geq 0$ :

$$\begin{array}{l} \max c \cdot x \\ \text{s.t. } Ax = b \\ \quad x \geq 0 \end{array}$$

It's fairly easy to convert any linear programming problem into slack form. Slack form is especially useful in executing the simplex algorithm (which we'll see in the next lecture).

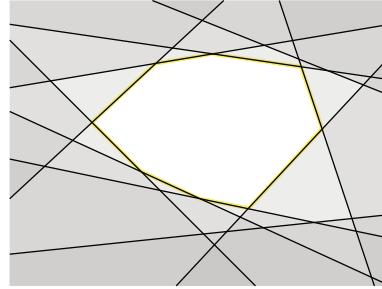
## 25.1 The Geometry of Linear Programming

A point  $x \in \mathbb{R}^d$  is *feasible* with respect to some linear programming problem if it satisfies all the linear constraints. The set of all feasible points is called the *feasible region* for that linear program. The feasible region has a particularly nice geometric structure that lends some useful intuition to the linear programming algorithms we'll see later.

<sup>1</sup>Confusingly, some authors call this *standard form*.

<sup>2</sup>Confusingly, some authors call this *standard form*.

Any linear equation in  $d$  variables defines a *hyperplane* in  $\mathbb{R}^d$ ; think of a line when  $d = 2$ , or a plane when  $d = 3$ . This hyperplane divides  $\mathbb{R}^d$  into two *halfspaces*; each halfspace is the set of points that satisfy some linear inequality. Thus, the set of feasible points is the intersection of several hyperplanes (one for each equality constraint) and halfspaces (one for each inequality constraint). The intersection of a finite number of hyperplanes and halfspaces is called a *polyhedron*. It's not hard to verify that any halfspace, and therefore any polyhedron, is *convex*—if a polyhedron contains two points  $x$  and  $y$ , then it contains the entire line segment  $\overline{xy}$ .



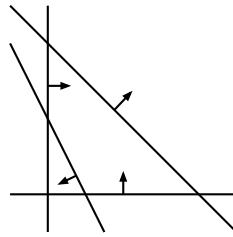
A two-dimensional polyhedron (white) defined by 10 linear inequalities.

By rotating  $\mathbb{R}^d$  (or choosing a coordinate frame) so that the objective function points downward, we can express *any* linear programming problem in the following geometric form:

Find the lowest point in a given polyhedron.

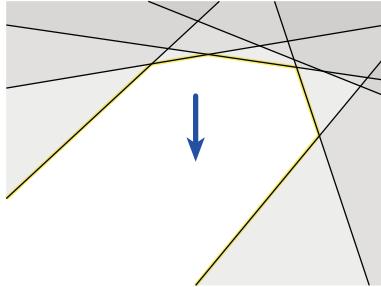
With this geometry in hand, we can easily picture two pathological cases where a given linear programming problem has no solution. The first possibility is that there are no feasible points; in this case the problem is called *infeasible*. For example, the following LP problem is infeasible:

$$\begin{aligned} & \text{maximize } x - y \\ & \text{subject to } 2x + y \leq 1 \\ & \quad x + y \geq 2 \\ & \quad x, y \geq 0 \end{aligned}$$



An infeasible linear programming problem; arrows indicate the constraints.

The second possibility is that there are feasible points at which the objective function is arbitrarily large; in this case, we call the problem *unbounded*. The same polyhedron could be unbounded for some objective functions but not others, or it could be unbounded for every objective function.



A two-dimensional polyhedron (white) that is unbounded downward but bounded upward.

## 25.2 Example 1: Shortest Paths

We can compute the length of the shortest path from  $s$  to  $t$  in a weighted directed graph by solving the following very simple linear programming problem.

$$\begin{aligned} & \text{maximize} && d_t \\ & \text{subject to} && d_s = 0 \\ & && d_v - d_u \leq \ell_{u \rightarrow v} \quad \text{for every edge } u \rightarrow v \end{aligned}$$

Here,  $\ell_{u \rightarrow v}$  is the length of the edge  $u \rightarrow v$ . Each variable  $d_v$  represents a tentative shortest-path distance from  $s$  to  $v$ . The constraints mirror the requirement that every edge in the graph must be relaxed. These relaxation constraints imply that in any feasible solution,  $d_v$  is *at most* the shortest path distance from  $s$  to  $v$ . Thus, somewhat counterintuitively, we are correctly *maximizing* the objective function to compute the *shortest* path! In the optimal solution, the objective function  $d_t$  is the actual shortest-path distance from  $s$  to  $t$ , but for any vertex  $v$  that is not on the shortest path from  $s$  to  $t$ ,  $d_v$  may be an underestimate of the true distance from  $s$  to  $v$ . However, we can obtain the true distances from  $s$  to every other vertex by modifying the objective function:

$$\begin{aligned} & \text{maximize} && \sum_v d_v \\ & \text{subject to} && d_s = 0 \\ & && d_v - d_u \leq \ell_{u \rightarrow v} \quad \text{for every edge } u \rightarrow v \end{aligned}$$

There is another formulation of shortest paths as an LP minimization problem using an indicator variable  $x_{u \rightarrow v}$  for each edge  $u \rightarrow v$ .

$$\begin{aligned} & \text{minimize} && \sum_{u \rightarrow v} \ell_{u \rightarrow v} \cdot x_{u \rightarrow v} \\ & \text{subject to} && \sum_u x_{u \rightarrow s} - \sum_w x_{s \rightarrow w} = 1 \\ & && \sum_u x_{u \rightarrow t} - \sum_w x_{t \rightarrow w} = -1 \\ & && \sum_u x_{u \rightarrow v} - \sum_w x_{v \rightarrow w} = 0 \quad \text{for every vertex } v \neq s, t \\ & && x_{u \rightarrow v} \geq 0 \quad \text{for every edge } u \rightarrow v \end{aligned}$$

Intuitively,  $x_{u \rightarrow v} = 1$  means  $u \rightarrow v$  lies on the shortest path from  $s$  to  $t$ , and  $x_{u \rightarrow v} = 0$  means  $u \rightarrow v$  does not lie on this shortest path. The constraints merely state that the path should start at  $s$ , end at  $t$ , and either pass through or avoid every other vertex  $v$ . Any path from  $s$  to  $t$ —in particular, the shortest path—clearly implies a feasible point for this linear program.

However, there are other feasible solutions, possibly even *optimal* solutions, with non-integral values that do not represent paths. Nevertheless, there is always an optimal solution in which every  $x_e$  is either 0 or 1 and the edges  $e$  with  $x_e = 1$  comprise the shortest path. (This fact is by no means obvious, but a proof is beyond the scope of these notes.) Moreover, in any optimal solution, even if not every  $x_e$  is an integer, the objective function gives the shortest path distance!

### 25.3 Example 2: Maximum Flows and Minimum Cuts

Recall that the input to the maximum  $(s, t)$ -flow problem consists of a weighted directed graph  $G = (V, E)$ , two special vertices  $s$  and  $t$ , and a function assigning a non-negative *capacity*  $c_e$  to each edge  $e$ . Our task is to choose the flow  $f_e$  across each edge  $e$ , as follows:

$$\begin{aligned} \text{maximize} \quad & \sum_w f_{s \rightarrow w} - \sum_u f_{u \rightarrow s} \\ \text{subject to} \quad & \sum_w f_{v \rightarrow w} - \sum_u f_{u \rightarrow v} = 0 \quad \text{for every vertex } v \neq s, t \\ & f_{u \rightarrow v} \leq c_{u \rightarrow v} \quad \text{for every edge } u \rightarrow v \\ & f_{u \rightarrow v} \geq 0 \quad \text{for every edge } u \rightarrow v \end{aligned}$$

Similarly, the minimum cut problem can be formulated using ‘indicator’ variables similarly to the shortest path problem. We have a variable  $S_v$  for each vertex  $v$ , indicating whether  $v \in S$  or  $v \in T$ , and a variable  $X_{u \rightarrow v}$  for each edge  $u \rightarrow v$ , indicating whether  $u \in S$  and  $v \in T$ , where  $(S, T)$  is some  $(s, t)$ -cut.<sup>3</sup>

$$\begin{aligned} \text{minimize} \quad & \sum_{u \rightarrow v} c_{u \rightarrow v} \cdot X_{u \rightarrow v} \\ \text{subject to} \quad & X_{u \rightarrow v} + S_v - S_u \geq 0 \quad \text{for every edge } u \rightarrow v \\ & X_{u \rightarrow v} \geq 0 \quad \text{for every edge } u \rightarrow v \\ & S_s = 1 \\ & S_t = 0 \end{aligned}$$

Like the minimization LP for shortest paths, there can be optimal solutions that assign fractional values to the variables. Nevertheless, the minimum value for the objective function is the cost of the minimum cut, and there is an optimal solution for which every variable is either 0 or 1, representing an actual minimum cut. No, this is not obvious; in particular, my claim is not a proof!

### 25.4 Linear Programming Duality

Each of these pairs of linear programming problems is related by a transformation called *duality*. For any linear programming problem, there is a corresponding dual linear program that can be obtained by a mechanical translation, essentially by swapping the constraints and the variables. The translation is simplest when the LP is in canonical form:

Primal (II)	Dual (II)
$\begin{aligned} \max \quad & c \cdot x \\ \text{s.t.} \quad & Ax \leq b \\ & x \geq 0 \end{aligned}$	$\begin{aligned} \min \quad & y \cdot b \\ \text{s.t.} \quad & yA \geq c \\ & y \geq 0 \end{aligned}$

<sup>3</sup>These two linear programs are not quite *syntactic* duals; I’ve added two redundant variables  $S_s$  and  $S_t$  to the min-cut program to increase readability.

We can also write the dual linear program in exactly the same canonical form as the primal, by swapping the coefficient vector  $c$  and the objective vector  $b$ , negating both vectors, and replacing the constraint matrix  $A$  with its negative transpose.<sup>4</sup>

Primal (II)	Dual (II)
$\begin{aligned} \max \quad & c \cdot x \\ \text{s.t. } & Ax \leq b \\ & x \geq 0 \end{aligned}$	$\begin{aligned} \max \quad & -b^\top \cdot y^\top \\ \text{s.t. } & -A^\top y^\top \leq -c \\ & y^\top \geq 0 \end{aligned}$

Written in this form, it should be immediately clear that duality is an *involution*: The dual of the dual linear program II is identical to the primal linear program II. The choice of which LP to call the ‘primal’ and which to call the ‘dual’ is totally arbitrary.<sup>5</sup>

**The Fundamental Theorem of Linear Programming.** *A linear program  $\Pi$  has an optimal solution  $x^*$  if and only if the dual linear program  $\Pi$  has an optimal solution  $y^*$  such that  $c \cdot x^* = y^* A x^* = y^* \cdot b$ .*

The weak form of this theorem is trivial to prove.

**Weak Duality Theorem.** *If  $x$  is a feasible solution for a canonical linear program  $\Pi$  and  $y$  is a feasible solution for its dual  $\Pi$ , then  $c \cdot x \leq y A x \leq y \cdot b$ .*

**Proof:** Because  $x$  is feasible for  $\Pi$ , we have  $Ax \leq b$ . Since  $y$  is positive, we can multiply both sides of the inequality to obtain  $yAx \leq y \cdot b$ . Conversely,  $y$  is feasible for  $\Pi$  and  $x$  is positive, so  $yAx \geq c \cdot x$ .  $\square$

It immediately follows that if  $c \cdot x = y \cdot b$ , then  $x$  and  $y$  are optimal solutions to their respective linear programs. This is in fact a fairly common way to prove that we have the optimal value for a linear program.

## 25.5 Duality Example

Before I prove the stronger duality theorem, let me first provide some intuition about where this duality thing comes from in the first place.<sup>6</sup> Consider the following linear programming problem:

$$\begin{aligned} \text{maximize} \quad & 4x_1 + x_2 + 3x_3 \\ \text{subject to} \quad & x_1 + 4x_2 \leq 2 \\ & 3x_1 - x_2 + x_3 \leq 4 \\ & x_1, x_2, x_3 \geq 0 \end{aligned}$$

Let  $\sigma^*$  denote the optimum objective value for this LP. The feasible solution  $x = (1, 0, 0)$  gives us a lower bound  $\sigma^* \geq 4$ . A different feasible solution  $x = (0, 0, 3)$  gives us a better lower bound  $\sigma^* \geq 9$ . We could play this game all day, finding different feasible solutions and getting ever larger lower bounds. How do we know when we’re done? Is there a way to prove an *upper bound* on  $\sigma^*$ ?

<sup>4</sup>For the notational purists: In these formulations,  $x$  and  $b$  are column vectors, and  $y$  and  $c$  are row vectors. This is a somewhat nonstandard choice. Yes, that means the dot in  $c \cdot x$  is redundant. Sue me.

<sup>5</sup>For historical reasons, maximization LPs tend to be called ‘primal’ and minimization LPs tend to be called ‘dual’. This is a pointless religious tradition, nothing more. Duality is a relationship between LP problems, not a type of LP problem.

<sup>6</sup>This example is taken from Robert Vanderbei’s excellent textbook *Linear Programming: Foundations and Extensions* [Springer, 2001], but the idea appears earlier in Jens Clausen’s 1997 paper ‘Teaching Duality in Linear Programming: The Multiplier Approach’.

In fact, there is. Let's multiply each of the constraints in our LP by a new non-negative scalar value  $y_i$ :

$$\begin{aligned} & \text{maximize} && 4x_1 + x_2 + 3x_3 \\ & \text{subject to} && y_1(x_1 + 4x_2) \leq 2y_1 \\ & && y_2(3x_1 - x_2 + x_3) \leq 4y_2 \\ & && x_1, x_2, x_3 \geq 0 \end{aligned}$$

Because each  $y_i$  is non-negative, we do not reverse any of the inequalities. Any feasible solution  $(x_1, x_2, x_3)$  must satisfy both of these inequalities, so it must also satisfy their sum:

$$(y_1 + 3y_2)x_1 + (4y_1 - y_2)x_2 + y_2x_3 \leq 2y_1 + 4y_2.$$

Now suppose that each  $y_i$  is larger than the  $i$ th coefficient of the objective function:

$$y_1 + 3y_2 \geq 4, \quad 4y_1 - y_2 \geq 1, \quad y_2 \geq 3.$$

This assumption lets us derive an upper bound on the objective value of *any* feasible solution:

$$4x_1 + x_2 + 3x_3 \leq (y_1 + 3y_2)x_1 + (4y_1 - y_2)x_2 + y_2x_3 \leq 2y_1 + 4y_2. \quad (*)$$

In particular, by plugging in the optimal solution  $(x_1^*, x_2^*, x_3^*)$  for the original LP, we obtain the following upper bound on  $\sigma^*$ :

$$\sigma^* = 4x_1^* + x_2^* + 3x_3^* \leq 2y_1 + 4y_2.$$

Now it's natural to ask how tight we can make this upper bound. How small can we make the expression  $2y_1 + 4y_2$  without violating any of the inequalities we used to prove the upper bound? This is just another linear programming problem.

$$\begin{aligned} & \text{minimize} && 2y_1 + 4y_2 \\ & \text{subject to} && y_1 + 3y_2 \geq 4 \\ & && 4y_1 - y_2 \geq 1 \\ & && y_2 \geq 3 \\ & && y_1, y_2 \geq 0 \end{aligned}$$

In fact, this is precisely the dual of our original linear program! Moreover, inequality  $(*)$  is just an instantiation of the Weak Duality Theorem.

## 25.6 Strong Duality

The Fundamental Theorem can be rephrased in the following form:

**Strong Duality Theorem.** *If  $x^*$  is an optimal solution for a canonical linear program  $\Pi$ , then there is an optimal solution  $y^*$  for its dual  $\Pi^*$ , such that  $c \cdot x^* = y^* A x^* = y^* \cdot b$ .*

**Proof (sketch):** I'll prove the theorem only for **non-degenerate** linear programs, in which (a) the optimal solution (if one exists) is a unique vertex of the feasible region, and (b) at most  $d$  constraint hyperplanes pass through any point. These non-degeneracy assumptions are relatively easy to enforce in practice and can be removed from the proof at the expense of some technical detail. I will also prove the theorem only for the case  $n \geq d$ ; the argument for under-constrained LPs is similar (if not simpler).

To develop some intuition, let's first consider the *very* special case where  $x^* = (0, 0, \dots, 0)$ . Let  $e_i$  denote the  $i$ th standard basis vector, whose  $i$ th coordinate is 1 and all other coordinates are 0. Because  $x_i^* = 0$  for all  $i$ , our non-degeneracy assumption implies the *strict* inequality  $a_i \cdot x^* < b_i$  for all  $i$ . Thus, any sufficiently small ball around the origin does not intersect any other constraint hyperplane  $a_i \cdot x = b_i$ . Thus, for all  $i$ , and for any sufficiently small  $\delta > 0$ , the vector  $\delta e_i$  is feasible. Because  $x^*$  is the unique optimum, we must have  $\delta c_i = c \cdot (\delta e_i) < c \cdot x^* = 0$ . We conclude that  $c_i < 0$  for all  $i$ .

Now let  $y = (0, 0, \dots, 0)$  as well. We immediately observe that  $yA \geq c$  and  $y \geq 0$ ; in other words,  $y$  is a *feasible* solution for the dual linear program  $\Pi$ . But  $y \cdot b = 0 = c \cdot x^*$ , so the weak duality theorem implies that  $y$  is an *optimal* solution to  $\Pi$ , and the proof is complete for this very special case.

Now let us consider the more general case. Let  $x^*$  be the optimal solution for the linear program  $\Pi$ ; our non-degeneracy assumption implies that this solution is unique, and that exactly  $d$  of the  $n$  linear constraints are satisfied with equality. Without loss of generality (by permuting the constraints and possibly changing coordinates), we can assume that these are the first  $d$  constraints. Thus, we have

$$\begin{aligned} a_i \cdot x^* &= b_i && \text{for all } i \leq d, \\ a_i \cdot x^* &< b_i && \text{for all } i \geq d+1, \end{aligned}$$

where  $a_i$  denotes the  $i$ th row of  $A$ . Let  $A_\bullet$  denote the  $d \times d$  matrix containing the first  $d$  rows of  $A$ . Our non-degeneracy assumption implies that  $A_\bullet$  has full rank, and thus has a well-defined inverse  $V = A_\bullet^{-1}$ .

Now define a vector  $y \in \mathbb{R}^n$  by setting

$$\begin{aligned} y_j &:= c \cdot v^j && \text{for all } j \leq d, \\ y_j &:= 0 && \text{for all } j \geq d+1, \end{aligned}$$

where  $v^j$  denotes the  $j$ th column of  $V = A_\bullet^{-1}$ . Note that  $a_i \cdot v^j = 0$  if  $i \neq j$ , and  $a_i \cdot v^j = 1$  if  $i = j$ .

To simplify notation, let  $y_\bullet = (y_1, y_2, \dots, y_d)$  and let  $b_\bullet = (b_1, b_2, \dots, b_d) = A_\bullet x^*$ . Because  $y_i = 0$  for all  $i \geq d+1$ , we immediately have

$$y \cdot b = y_\bullet \cdot b_\bullet = c V b_\bullet = c A_\bullet^{-1} b_\bullet = c \cdot x^*$$

and

$$yA = y_\bullet A_\bullet = c V A_\bullet = c A_\bullet^{-1} A_\bullet = c.$$

The point  $x^*$  lies on exactly  $d$  constraint hyperplanes; moreover, any sufficiently small ball around  $x^*$  intersects *only* those  $d$  constraint hyperplanes. Consider the point  $\tilde{x} = x^* - \varepsilon v^j$ , for some index  $1 \leq j \leq d$  and some sufficiently small  $\varepsilon > 0$ . We have  $a_i \cdot \tilde{x} = a_i \cdot x^* - \varepsilon(a_i \cdot v^j) = b_i$  for all  $i \neq j$ , and  $a_j \cdot \tilde{x} = a_j \cdot x^* - \varepsilon(a_j \cdot v^j) = b_j - \varepsilon < b_j$ . Thus,  $\tilde{x}$  is a feasible point for  $\Pi$ . Because  $x^*$  is the *unique* optimum for  $\Pi$ , we must have  $c \cdot \tilde{x} = c \cdot x^* - \varepsilon(c \cdot v^j) < c \cdot x^*$ . We conclude that  $y_j = c \cdot v^j > 0$  for all  $j$ .

We have shown that  $yA \geq c$  and  $y \geq 0$ , so  $y$  is a *feasible* solution for the dual linear program  $\Pi$ . We have also shown that  $y \cdot b = c \cdot x^*$ , so by the Weak Duality Theorem,  $y$  is also an *optimal* solution for  $\Pi$ , and the proof is complete!  $\square$

We can also give a useful geometric interpretation to the vector  $y_\bullet \in \mathbb{R}^d$ . Each linear equation  $a_i \cdot x = b_i$  defines a hyperplane in  $\mathbb{R}^d$  with normal vector  $a_i$ . The normal vectors  $a_1, \dots, a_d$  are linearly independent (by non-degeneracy) and therefore describe a coordinate frame for the vector space  $\mathbb{R}^d$ . The definition of  $y_\bullet$  implies that  $c = y_\bullet A_\bullet = \sum_{i=1}^d y_i a_i$ . In other words,  $y_\bullet$  lists the coefficients of the objective vector  $c$  in the coordinate frame  $a_1, \dots, a_d$ .

## 25.7 Complementary Slackness

**Complementary Slackness Theorem.** Let  $x^*$  be an optimal solution to a canonical linear program  $\Pi$ , and let  $y^*$  be an optimal solution to its dual  $\Pi$ . Then for every index  $i$ , we have  $y_i^* > 0$  if and only if  $a_i \cdot x^* = b_i$ . Symmetrically, for every index  $j$ , we have  $x_j^* > 0$  if and only if  $y^* \cdot a^j = c_j$ .

To be written

## Exercises

1. (a) Describe how to transform any linear program written in general form into an equivalent linear program written in slack form.

$\begin{aligned} & \text{maximize} && \sum_{j=1}^d c_j x_j \\ & \text{subject to} && \sum_{j=1}^d a_{ij} x_j \leq b_i \quad \text{for each } i = 1 \dots p \\ & && \sum_{j=1}^d a_{ij} x_j = b_i \quad \text{for each } i = p+1 \dots p+q \\ & && \sum_{j=1}^d a_{ij} x_j \geq b_i \quad \text{for each } i = p+q+1 \dots n \end{aligned}$	$\implies$ <div style="border: 1px solid black; padding: 5px; display: inline-block;"> <math display="block">\begin{aligned} &amp; \max &amp;&amp; c \cdot x \\ &amp; \text{s.t.} &amp;&amp; Ax = b \\ &amp; &amp;&amp; x \geq 0 \end{aligned}</math> </div>
--	--

- (b) Describe precisely how to dualize a linear program written in slack form.

- (c) Describe precisely how to dualize a linear program written in general form:

In all cases, keep the number of variables in the resulting linear program as small as possible.

2. A matrix  $A = (a_{ij})$  is *skew-symmetric* if and only if  $a_{ji} = -a_{ij}$  for all indices  $i \neq j$ ; in particular, every skew-symmetric matrix is square. A canonical linear program  $\max\{c \cdot x \mid Ax \leq b; x \geq 0\}$  is *self-dual* if the matrix  $A$  is skew-symmetric and the objective vector  $c$  is *equal* to the constraint vector  $b$ .
  - (a) Prove that any self-dual linear program  $\Pi$  is syntactically equivalent to its dual program  $\Pi$ .
  - (b) Show that any linear program  $\Pi$  with  $d$  variables and  $n$  constraints can be transformed into a self-dual linear program with  $n+d$  variables and  $n+d$  constraints. The optimal solution to the self-dual program should include both the optimal solution for  $\Pi$  (in  $d$  of the variables) and the optimal solution for the dual program  $\Pi$  (in the other  $n$  variables).
3. (a) Give a linear-programming formulation of the *maximum-cardinality bipartite matching* problem. The input is a bipartite graph  $G = (U \cup V; E)$ , where  $E \subseteq U \times V$ ; the output is the largest matching in  $G$ . Your linear program should have one variable for each edge.
  - (b) Now dualize the linear program from part (a). What do the dual variables represent? What does the objective function represent? What problem is this!?
4. Give a linear-programming formulation of the *minimum-cost feasible circulation problem*. Here you are given a flow network whose edges have both capacities and costs, and your goal is to find a feasible circulation (flow with value 0) whose cost is as small as possible.

5. An *integer program* is a linear program with the additional constraint that the variables must take only integer values.

- (a) Prove that deciding whether an integer program has a feasible solution is NP-complete.
- (b) Prove that finding the optimal feasible solution to an integer program is NP-hard.

[Hint: Almost any NP-hard decision problem can be formulated as an integer program. Pick your favorite.]

\*6. Helly's theorem states that for any collection of convex bodies in  $\mathbb{R}^d$ , if every  $d + 1$  of them intersect, then there is a point lying in the intersection of all of them. Prove Helly's theorem for the special case where the convex bodies are halfspaces. Equivalently, show that if a system of linear inequalities  $Ax \leq b$  does not have a solution, then we can select  $d + 1$  of the inequalities such that the resulting subsystem also does not have a solution. [Hint: Construct a dual LP from the system by choosing a 0 cost vector.]

7. Given points  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  in the plane, the *linear regression problem* asks for real numbers  $a$  and  $b$  such that the line  $y = ax + b$  fits the points as closely as possible, according to some criterion. The most common fit criterion is minimizing the  $L_2$  error, defined as follows:<sup>7</sup>

$$\varepsilon_2(a, b) = \sum_{i=1}^n (y_i - ax_i - b)^2.$$

But there are several other fit criteria, some of which can be optimized via linear programming.

- (a) The  $L_1$  error (or *total absolute deviation*) of the line  $y = ax + b$  is defined as follows:

$$\varepsilon_1(a, b) = \sum_{i=1}^n |y_i - ax_i - b|.$$

Describe a linear program whose solution  $(a, b)$  describes the line with minimum  $L_1$  error.

- (b) The  $L_\infty$  error (or *maximum absolute deviation*) of the line  $y = ax + b$  is defined as follows:

$$\varepsilon_\infty(a, b) = \max_{i=1}^n |y_i - ax_i - b|.$$

Describe a linear program whose solution  $(a, b)$  describes the line with minimum  $L_\infty$  error.

---

<sup>7</sup>This measure is also known as *sum of squared residuals*, and the algorithm to compute the best fit is normally called (*ordinary/linear*) *least squares fitting*.

*Simplicibus itaque verbis gaudet Mathematica Veritas, cum etiam per se simplex sit Veritatis oratio. [And thus Mathematical Truth prefers simple words, because the language of Truth is itself simple.]*

— Tycho Brahe (quoting Seneca (quoting Euripides))  
*Epistolarum astronomicarum liber primus* (1596)

*When a jar is broken, the space that was inside  
 Merges into the space outside.*

*In the same way, my mind has merged in God;  
 To me, there appears no duality.*

— Sankara, *Viveka-Chudamani* (c. 700), translator unknown

## \*26 Linear Programming Algorithms

In this lecture, we'll see a few algorithms for actually solving linear programming problems. The most famous of these, the *simplex method*, was proposed by George Dantzig in 1947. Although most variants of the simplex algorithm performs well in practice, no deterministic simplex variant is known to run in sub-exponential time in the worst case.<sup>1</sup> However, if the dimension of the problem is considered a constant, there are several linear programming algorithms that run in *linear* time. I'll describe a particularly simple randomized algorithm due to Raimund Seidel.

My approach to describing these algorithms will rely much more heavily on geometric intuition than the usual linear-algebraic formalism. This works better for me, but your mileage may vary. For a more traditional description of the simplex algorithm, see Robert Vanderbei's excellent textbook *Linear Programming: Foundations and Extensions* [Springer, 2001], which can be freely downloaded (but not legally printed) from the author's website.

### 26.1 Bases, Feasibility, and Local Optimality

Consider the canonical linear program  $\max\{c \cdot x \mid Ax \leq b, x \geq 0\}$ , where  $A$  is an  $n \times d$  constraint matrix,  $b$  is an  $n$ -dimensional coefficient vector, and  $c$  is a  $d$ -dimensional objective vector. We will interpret this linear program geometrically as looking for the lowest point in a convex polyhedron in  $\mathbb{R}^d$ , described as the intersection of  $n + d$  halfspaces. As in the last lecture, we will consider only *non-degenerate* linear programs: Every subset of  $d$  constraint hyperplanes intersect in a single point; at most  $d$  constraint hyperplanes pass through any point; and objective vector is linearly independent from any  $d - 1$  constraint vectors.

A **basis** is a subset of  $d$  constraints, which by our non-degeneracy assumption must be linearly independent. The **location** of a basis is the unique point  $x$  that satisfies all  $d$  constraints with equality; geometrically,  $x$  is the unique intersection point of the  $d$  hyperplanes. The **value** of a basis is  $c \cdot x$ , where  $x$  is the location of the basis. There are precisely  $\binom{n+d}{d}$  bases. Geometrically, the set of constraint hyperplanes defines a decomposition of  $\mathbb{R}^d$  into convex polyhedra; this cell decomposition is called the **arrangement** of the hyperplanes. Every subset of  $d$  hyperplanes (that is, every basis) defines a **vertex** of this arrangement (the location of the basis). I will use the words 'vertex' and 'basis' interchangeably.

A basis is **feasible** if its location  $x$  satisfies all the linear constraints, or geometrically, if the point  $x$  is a vertex of the polyhedron. If there are no feasible bases, the linear program is **infeasible**.

<sup>1</sup>However, there are *randomized* variants of the simplex algorithm that run in subexponential *expected* time, most notably the RANDOMFACET algorithm analyzed by Gil Kalai in 1992, and independently by Jiří Matoušek, Micha Sharir, and Emo Welzl in 1996. No randomized variant is known to run in polynomial time. In particular, in 2010, Oliver Friedmann, Thomas Dueholm Hansen, and Uri Zwick proved that the worst-case expected running time of RANDOMFACET is superpolynomial.

A basis is *locally optimal* if its location  $x$  is the optimal solution to the linear program with the same objective function and *only* the constraints in the basis. Geometrically, a basis is locally optimal if its location  $x$  is the lowest point in the intersection of those  $d$  halfspaces. A careful reading of the proof of the Strong Duality Theorem reveals that local optimality is the dual equivalent of feasibility; a basis is locally feasible for a linear program  $\Pi$  if and only if the same basis is feasible for the dual linear program  $\Pi^*$ . For this reason, locally optimal bases are sometimes also called *dual feasible*. If there are no locally optimal bases, the linear program is *unbounded*.<sup>2</sup>

Two bases are *neighbors* if they have  $d - 1$  constraints in common. Equivalently, in geometric terms, two vertices are neighbors if they lie on a *line* determined by some  $d - 1$  constraint hyperplanes. Every basis is a neighbor of exactly  $d n$  other bases; to change a basis into one of its neighbors, there are  $d$  choices for which constraint to remove and  $n$  choices for which constraint to add. The graph of vertices and edges on the boundary of the feasible polyhedron is a subgraph of the basis graph.

The Weak Duality Theorem implies that the value of every feasible basis is less than or equal to the value of every locally optimal basis; equivalently, every feasible vertex is higher than every locally optimal vertex. The Strong Duality Theorem implies that (under our non-degeneracy assumption), if a linear program has an optimal solution, it is the *unique* vertex that is both feasible and locally optimal. Moreover, the optimal solution is both the lowest feasible vertex and the highest locally optimal vertex.

## 26.2 The Primal Simplex Algorithm: Falling Marbles

From a geometric standpoint, Dantzig's simplex algorithm is very simple. The input is a set  $H$  of halfspaces; we want the lowest vertex in the intersection of these halfspaces.

```
SIMPLEX1( $H$ ):
  if  $\cap H = \emptyset$ 
    return INFEASIBLE
   $x \leftarrow$  any feasible vertex
  while  $x$  is not locally optimal
    <pivot downward, maintaining feasibility>
    if every feasible neighbor of  $x$  is higher than  $x$ 
      return UNBOUNDED
    else
       $x \leftarrow$  any feasible neighbor of  $x$  that is lower than  $x$ 
  return  $x$ 
```

Let's ignore the first three lines for the moment. The algorithm maintains a feasible vertex  $x$ . At each so-called *pivot* operation, the algorithm moves to a *lower* vertex, so the algorithm never visits the same vertex more than once. Thus, the algorithm must halt after at most  $\binom{n+d}{d}$  pivots. When the algorithm halts, either the feasible vertex  $x$  is locally optimal, and therefore the optimum vertex, or the feasible vertex  $x$  is not locally optimal but has no lower feasible neighbor, in which case the feasible region must be unbounded.

Notice that we have not specified *which* neighbor to choose at each pivot. Many different pivoting rules have been proposed, but for almost every known pivot rule, there is an input polyhedron that requires an exponential number of pivots under that rule. No pivoting rule is known that guarantees a polynomial number of pivots in the worst case, or even in expectation.<sup>3</sup>

<sup>2</sup>For non-degenerate linear programs, the feasible region is unbounded in the objective direction if and only if no basis is locally optimal. However, there are degenerate linear programs with no locally optimal basis that are infeasible.

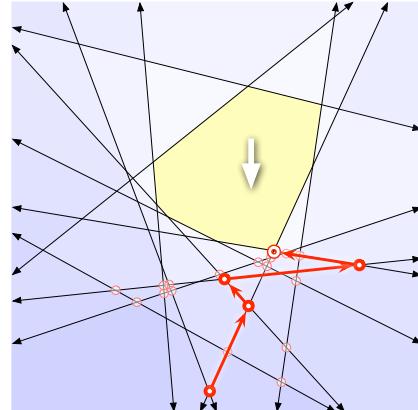
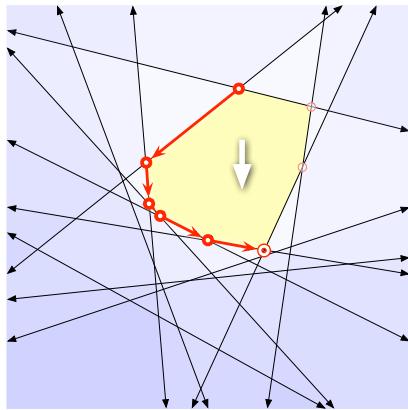
<sup>3</sup>In 1957, Hirsch conjectured that for *any* linear programming instance with  $d$  variables and  $n + d$  constraints, starting at any feasible basis, there is a sequence of *at most*  $n$  pivots that leads to the optimal basis. This long-standing conjecture was finally disproved in 2010 by Francisco Santos, who described an counterexample with 43 variables, 86 facets, and diameter 44.

### 26.3 The Dual Simplex Algorithm: Rising Bubbles

We can also geometrically interpret the execution of the simplex algorithm on the dual linear program II. Again, the input is a set  $H$  of halfspaces, and we want the lowest vertex in the intersection of these halfspaces. By the Strong Duality Theorem, this is the same as the *highest locally-optimal* vertex in the hyperplane arrangement.

```
SIMPLEX2( $H$ ):
    if there is no locally optimal vertex
        return UNBOUNDED
     $x \leftarrow$  any locally optimal vertex
    while  $x$  is not feasible
        «pivot upward, maintaining local optimality»
        if every locally optimal neighbor of  $x$  is lower than  $x$ 
            return INFEASIBLE
        else
             $x \leftarrow$  any locally-optimal neighbor of  $x$  that is higher than  $x$ 
    return  $x$ 
```

Let's ignore the first three lines for the moment. The algorithm maintains a locally optimal vertex  $x$ . At each pivot operation, it moves to a *higher* vertex, so the algorithm never visits the same vertex more than once. Thus, the algorithm must halt after at most  $\binom{n+d}{d}$  pivots. When the algorithm halts, either the locally optimal vertex  $x$  is feasible, and therefore the optimum vertex, or the locally optimal vertex  $x$  is not feasible but has no higher locally optimal neighbor, in which case the problem must be infeasible.



The primal simplex (falling marble) algorithm in action.    The dual simplex (rising bubble) algorithm in action.

From the standpoint of linear algebra, there is absolutely no difference between running **SIMPLEX1** on any linear program  $\Pi$  and running **SIMPLEX2** on the dual linear program  $\Pi^*$ . The actual code is identical. The only difference between the two algorithms is how we interpret the linear algebra geometrically.

### 26.4 Computing the Initial Basis

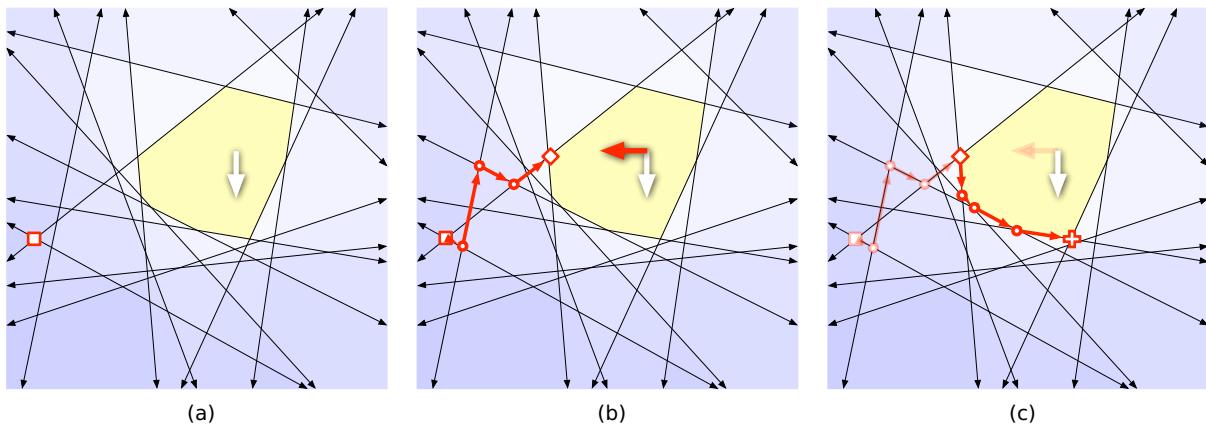
To complete our description of the simplex algorithm, we need to describe how to find the initial vertex  $x$ . Our algorithm relies on the following simple observations.

First, the feasibility of a vertex does not depend at all on the choice of objective vector; a vertex is either feasible for every objective function or for none. No matter how we rotate the polyhedron, every feasible vertex stays feasible. Conversely (or by duality, equivalently), the local optimality of a vertex

does not depend on the exact location of the  $d$  hyperplanes, but only on their normal directions and the objective function. No matter how we translate the hyperplanes, every locally optimal vertex stays locally optimal. In terms of the original matrix formulation, feasibility depends on  $A$  and  $b$  but not  $c$ , and local optimality depends on  $A$  and  $c$  but not  $b$ .

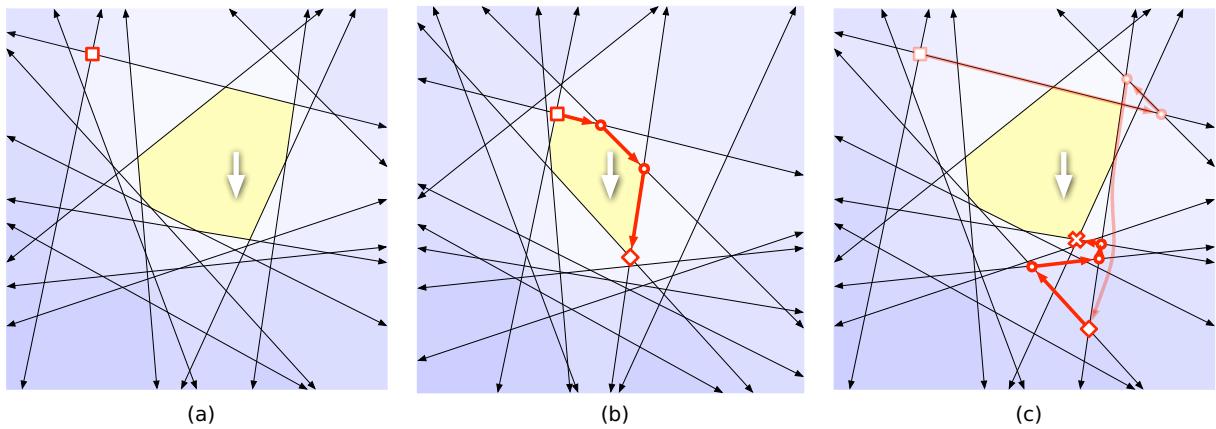
The second important observation is that *every* basis is locally optimal for *some* objective function. Specifically, it suffices to choose any vector that has a positive inner product with each of the normal vectors of the  $d$  chosen hyperplanes. Equivalently, we can make *any* basis feasible by translating the hyperplanes appropriately. Specifically, it suffices to translate the chosen  $d$  hyperplanes so that they pass through the origin, and then translate all the other halfspaces so that they strictly contain the origin.

Our strategy for finding our initial feasible vertex is to choose *any* vertex, choose a new objective function that makes that vertex locally optimal, and then find the optimal vertex for *that* objective function by running the (dual) simplex algorithm. This vertex must be feasible, even after we restore the original objective function!



(a) Choose any basis. (b) Rotate objective to make it locally optimal, and pivot 'upward' to find a feasible basis.  
 (c) Pivot downward to the optimum basis for the original objective.

Equivalently, to find an initial locally optimal vertex, we choose *any* vertex, translate the hyperplanes so that that vertex becomes feasible, and then find the optimal vertex for those translated constraints using the (primal) simplex algorithm. This vertex must be locally optimal, even after we restore the hyperplanes to their original locations!



(a) Choose any basis. (b) Translate constraints to make it feasible, and pivot downward to find a locally optimal basis.  
 (c) Pivot upward to the optimum basis for the original constraints.

Here are more complete descriptions of the simplex algorithm with this initialization rule, in both primal and dual forms. As usual, the input is a set  $H$  of halfspaces, and the algorithms either return the lowest vertex in the intersection of these halfspaces or report that no such vertex exists.

```
SIMPLEX1( $H$ ):
 $x \leftarrow$  any vertex
 $\tilde{H} \leftarrow$  any rotation of  $H$  that makes  $x$  locally optimal
while  $x$  is not feasible
    if every locally optimal neighbor of  $x$  is lower (wrt  $\tilde{H}$ ) than  $x$ 
        return INFEASIBLE
    else
         $x \leftarrow$  any locally optimal neighbor of  $x$  that is higher (wrt  $\tilde{H}$ ) than  $x$ 
while  $x$  is not locally optimal
    if every feasible neighbor of  $x$  is higher than  $x$ 
        return UNBOUNDED
    else
         $x \leftarrow$  any feasible neighbor of  $x$  that is lower than  $x$ 
return  $x$ 
```

```
SIMPLEX2( $H$ ):
 $x \leftarrow$  any vertex
 $\tilde{H} \leftarrow$  any translation of  $H$  that makes  $x$  feasible
while  $x$  is not locally optimal
    if every feasible neighbor of  $x$  is higher (wrt  $\tilde{H}$ ) than  $x$ 
        return UNBOUNDED
    else
         $x \leftarrow$  any feasible neighbor of  $x$  that is lower (wrt  $\tilde{H}$ ) than  $x$ 
while  $x$  is not feasible
    if every locally optimal neighbor of  $x$  is lower than  $x$ 
        return INFEASIBLE
    else
         $x \leftarrow$  any locally-optimal neighbor of  $x$  that is higher than  $x$ 
return  $x$ 
```

## 26.5 Linear Expected Time for Fixed Dimensions

In most geometric applications of linear programming, the number of variables is a small constant, but the number of constraints may still be very large.

The input to the following algorithm is a set  $H$  of  $n$  halfspaces and a set  $B$  of  $b$  hyperplanes. ( $B$  stands for *basis*.) The algorithm returns the lowest point in the intersection of the halfspaces in  $H$  and the hyperplanes  $B$ . At the top level of recursion,  $B$  is empty. I will implicitly assume that the linear program is both feasible and bounded. (If necessary, we can guarantee boundedness by adding a single halfspace to  $H$ , and we can guarantee feasibility by adding a dimension.) A point  $x$  **violates** a constraint  $h$  if it is not contained in the corresponding halfspace.

```

SEIDELLP( $H, B$ ) :
  if  $|B| = d$ 
    return  $\bigcap B$ 
  if  $|H \cup B| = d$ 
    return  $\bigcap(H \cup B)$ 
   $h \leftarrow$  random element of  $H$ 
   $x \leftarrow$  SEIDELLP( $H \setminus h, B$ )      (*)
  if  $x$  violates  $h$ 
    return SEIDELLP( $H \setminus h, B \cup \partial h$ )
  else
    return  $x$ 

```

The point  $x$  recursively computed in line (\*) is the optimal solution if and only if the random halfspace  $h$  is *not* one of the  $d$  halfspaces that define the optimal solution. In other words, the probability of calling  $\text{SEIDELLP}(H, B \cup h)$  is exactly  $(d - b)/n$ . Thus, we have the following recurrence for the expected number of recursive calls for this algorithm:

$$T(n, b) = \begin{cases} 1 & \text{if } b = d \text{ or } n + b = d \\ T(n - 1, b) + \frac{d - b}{n} \cdot T(n - 1, b + 1) & \text{otherwise} \end{cases}$$

The recurrence is somewhat simpler if we write  $\delta = d - b$ :

$$T(n, \delta) = \begin{cases} 1 & \text{if } \delta = 0 \text{ or } n = \delta \\ T(n - 1, \delta) + \frac{\delta}{n} \cdot T(n - 1, \delta - 1) & \text{otherwise} \end{cases}$$

It's easy to prove by induction that  $T(n, \delta) = O(\delta! n)$ :

$$\begin{aligned} T(n, \delta) &= T(n - 1, \delta) + \frac{\delta}{n} \cdot T(n - 1, \delta - 1) \\ &\leq \delta!(n - 1) + \frac{\delta}{n}(\delta - 1)! \cdot (n - 1) && [\text{induction hypothesis}] \\ &= \delta!(n - 1) + \delta! \frac{n - 1}{n} \\ &\leq \delta! n \end{aligned}$$

At the top level of recursion, we perform one violation test in  $O(d)$  time. In each of the base cases, we spend  $O(d^3)$  time computing the intersection point of  $d$  hyperplanes, and in the first base case, we spend  $O(dn)$  additional time testing for violations. More careful analysis implies that the algorithm runs in  $O(d! \cdot n)$  ***expected time***.

## Exercises

1. Fix a non-degenerate linear program in canonical form with  $d$  variables and  $n + d$  constraints.
  - (a) Prove that every *feasible* basis has exactly  $d$  *feasible* neighbors.
  - (b) Prove that every *locally optimal* basis has exactly  $n$  *locally optimal* neighbors.

2. Suppose you have a subroutine that can solve linear programs in polynomial time, but only if they are both feasible and bounded. Describe an algorithm that solves *arbitrary* linear programs in polynomial time. Your algorithm should return an optimal solution if one exists; if no optimum exists, your algorithm should report that the input instance is UNBOUNDED or INFEASIBLE, whichever is appropriate. [*Hint: Add one variable and one constraint.*]
3. (a) Give an example of a non-empty polyhedron  $Ax \leq b$  that is unbounded for every objective vector  $c$ .  
(b) Give an example of an infeasible linear program whose dual is also infeasible.  
In both cases, your linear program will be degenerate.
4. Describe and analyze an algorithm that solves the following problem in  $O(n)$  time: Given  $n$  red points and  $n$  blue points in the plane, either find a line that separates every red point from every blue point, or prove that no such line exists.
5. The single-source shortest path problem can be formulated as a linear programming problem, with one variable  $d_v$  for each vertex  $v \neq s$  in the input graph, as follows:

$$\begin{aligned} &\text{maximize} \quad \sum_v d_v \\ &\text{subject to} \quad d_v \leq \ell_{s \rightarrow v} \quad \text{for every edge } s \rightarrow v \\ &\quad d_v - d_u \leq \ell_{u \rightarrow v} \quad \text{for every edge } u \rightarrow v \text{ with } u \neq s \\ &\quad d_v \geq 0 \quad \text{for every vertex } v \neq s \end{aligned}$$

This problem asks you to describe the behavior of the simplex algorithm on this linear program in terms of distances. Assume that the edge weights  $\ell_{u \rightarrow v}$  are all non-negative and that there is a unique shortest path between any two vertices in the graph.

- (a) What is a basis for this linear program? What is a feasible basis? What is a locally optimal basis?  
(b) Show that in the optimal basis, every variable  $d_v$  is equal to the shortest-path distance from  $s$  to  $v$ .  
(c) Describe the primal simplex algorithm for the shortest-path linear program directly in terms of vertex distances. In particular, what does it mean to pivot from a feasible basis to a neighboring feasible basis, and how can we execute such a pivot quickly?  
(d) Describe the dual simplex algorithm for the shortest-path linear program directly in terms of vertex distances. In particular, what does it mean to pivot from a locally optimal basis to a neighboring locally optimal basis, and how can we execute such a pivot quickly?  
(e) Is Dijkstra's algorithm an instance of the simplex method? Justify your answer.  
(f) Is Shimbel's algorithm an instance of the simplex method? Justify your answer.
6. The maximum  $(s, t)$ -flow problem can be formulated as a linear programming problem, with one

variable  $f_{u \rightarrow v}$  for each edge  $u \rightarrow v$  in the input graph:

$$\begin{aligned} & \text{maximize} \quad \sum_w f_{s \rightarrow w} - \sum_u f_{u \rightarrow s} \\ & \text{subject to} \quad \sum_w f_{v \rightarrow w} - \sum_u f_{u \rightarrow v} = 0 \quad \text{for every vertex } v \neq s, t \\ & \quad \quad \quad f_{u \rightarrow v} \leq c_{u \rightarrow v} \quad \text{for every edge } u \rightarrow v \\ & \quad \quad \quad f_{u \rightarrow v} \geq 0 \quad \text{for every edge } u \rightarrow v \end{aligned}$$

This problem asks you to describe the behavior of the simplex algorithm on this linear program in terms of flows.

- (a) What is a basis for this linear program? What is a feasible basis? What is a locally optimal basis?
  - (b) Show that the optimal basis represents a maximum flow.
  - (c) Describe the primal simplex algorithm for the flow linear program directly in terms of flows. In particular, what does it mean to pivot from a feasible basis to a neighboring feasible basis, and how can we execute such a pivot quickly?
  - (d) Describe the dual simplex algorithm for the flow linear program directly in terms of flows. In particular, what does it mean to pivot from a locally optimal basis to a neighboring locally optimal basis, and how can we execute such a pivot quickly?
  - (e) Is the Ford-Fulkerson augmenting path algorithm an instance of the simplex method? Justify your answer. *[Hint: There is a one-line argument.]*
7. (a) Formulate the minimum spanning tree problem as an instance of linear programming. Try to minimize the number of variables and constraints.
- (b) In your MST linear program, what is a basis? What is a feasible basis? What is a locally optimal basis?
- (c) Describe the primal simplex algorithm for your MST linear program directly in terms of the input graph. In particular, what does it mean to pivot from a feasible basis to a neighboring feasible basis, and how can we execute such a pivot quickly?
- (d) Describe the dual simplex algorithm for your MST linear program directly in terms of the input graph. In particular, what does it mean to pivot from a locally optimal basis to a neighboring locally optimal basis, and how can we execute such a pivot quickly?
- (e) Which of the classical MST algorithms (Borůvka, Jarník, Kruskal, reverse greedy), if any, are instances of the simplex method? Justify your answer.

*It was a Game called Yes and No, where Scrooge's nephew had to think of something, and the rest must find out what; he only answering to their questions yes or no, as the case was. The brisk fire of questioning to which he was exposed, elicited from him that he was thinking of an animal, a live animal, rather a disagreeable animal, a savage animal, an animal that growled and grunted sometimes, and talked sometimes, and lived in London, and walked about the streets, and wasn't made a show of, and wasn't led by anybody, and didn't live in a menagerie, and was never killed in a market, and was not a horse, or an ass, or a cow, or a bull, or a tiger, or a dog, or a pig, or a cat, or a bear. At every fresh question that was put to him, this nephew burst into a fresh roar of laughter; and was so inexpressibly tickled, that he was obliged to get up off the sofa and stamp. At last the plump sister, falling into a similar state, cried out :*

*"I have found it out! I know what it is, Fred! I know what it is!"*

*"What is it?" cried Fred.*

*"It's your Uncle Scro-o-o-o-oge!"*

*Which it certainly was. Admiration was the universal sentiment, though some objected that the reply to "Is it a bear?" ought to have been "Yes;" inasmuch as an answer in the negative was sufficient to have diverted their thoughts from Mr Scrooge, supposing they had ever had any tendency that way.*

— Charles Dickens, *A Christmas Carol* (1843)

## 27 Lower Bounds

### 27.1 Huh? Whuzzat?

So far in this class we've been developing algorithms and data structures to solve certain problems as quickly as possible. Starting with this lecture, we'll turn the tables, by proving that certain problems *cannot* be solved as quickly as we might like them to be.

Let  $T_A(X)$  denote the running time of algorithm  $A$  given input  $X$ . For most of the semester, we've been concerned with the worst-case running time of  $A$  as a function of the input size:

$$T_A(n) := \max_{|X|=n} T_A(X).$$

The worst-case complexity of a problem  $\Pi$  is the worst-case running time of the *fastest* algorithm for solving it:

$$T_\Pi(n) := \min_{A \text{ solves } \Pi} T_A(n) = \min_{A \text{ solves } \Pi} \max_{|X|=n} T_A(X).$$

Any algorithm  $A$  that solves  $\Pi$  immediately implies an *upper bound* on the complexity of  $\Pi$ ; the inequality  $T_\Pi(n) \leq T_A(n)$  follows directly from the definition of  $T_\Pi$ . Just as obviously, faster algorithms give us better (smaller) upper bounds. In other words, whenever we give a running time for an algorithm, what we're really doing—and what most computer scientists devote their entire careers doing<sup>1</sup>—is bragging about how *easy* some problem is.

Now, instead of bragging about how easy problems are, we will argue that certain problems are *hard*, by proving *lower bounds* on their complexity. This is considerably harder than proving an upper bound, because it's no longer enough to examine a single algorithm. To prove an inequality of the form

---

<sup>1</sup>This sometimes leads to long sequences of results that sound like an obscure version of “Name that Tune”:

Lennes: “I can triangulate that polygon in  $O(n^2)$  time.”

Shamos: “I can triangulate that polygon in  $O(n \log n)$  time.”

Tarjan: “I can triangulate that polygon in  $O(n \log \log n)$  time.”

Seidel: “I can triangulate that polygon in  $O(n \log^* n)$  time.” [Audience gasps.]

Chazelle: “I can triangulate that polygon in  $O(n)$  time.” [Audience gasps and applauds.]

“Triangulate that polygon!”

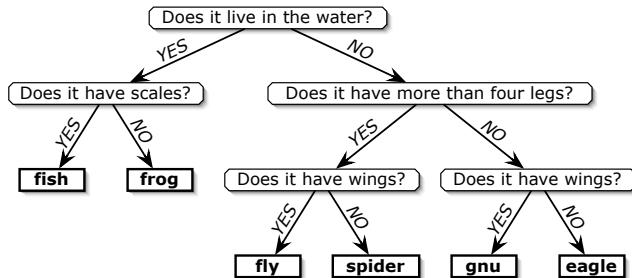
$T_{\Pi}(n) = \Omega(f(n))$ , we must prove that *every* algorithm that solves  $\Pi$  has a worst-case running time  $\Omega(f(n))$ , or equivalently, that *no* algorithm runs in  $o(f(n))$  time.

## 27.2 Decision Trees

Unfortunately, there is no formal definition of the phrase ‘all algorithms’!<sup>2</sup> So when we derive lower bounds, we first have to specify *precisely* what kinds of algorithms we will consider and *precisely* how to measure their running time. This specification is called a *model of computation*.

One rather powerful model of computation—and the only model we’ll talk about in this lecture—is the *decision tree* model. A decision tree is, as the name suggests, a tree. Each internal node in the tree is labeled by a *query*, which is just a question about the input. The edges out of a node correspond to the possible answers to that node’s query. Each leaf of the tree is labeled with an *output*. To compute with a decision tree, we start at the root and follow a path down to a leaf. At each internal node, the answer to the query tells us which node to visit next. When we reach a leaf, we output its label.

For example, the guessing game where one person thinks of an animal and the other person tries to figure it out with a series of yes/no questions can be modeled as a decision tree. Each internal node is labeled with a question and has two edges labeled ‘yes’ and ‘no’. Each leaf is labeled with an animal.

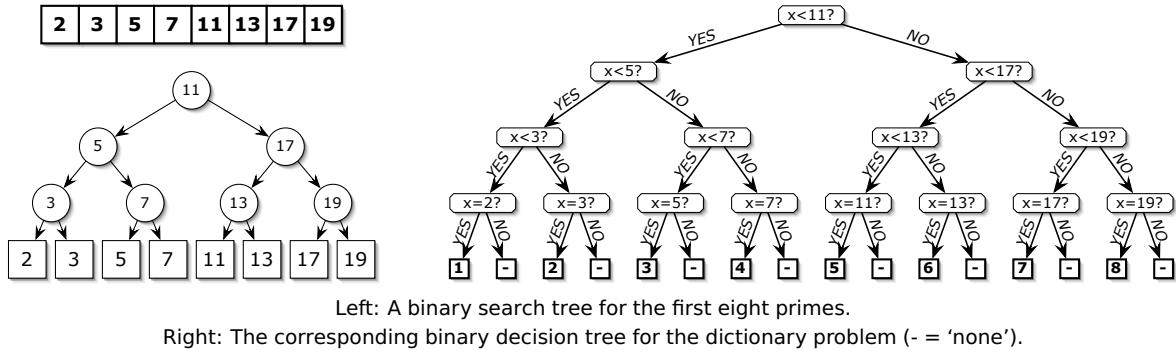


A decision tree to choose one of six animals.

Here’s another simple and familiar example, called the *dictionary problem*. Let  $A$  be a fixed array with  $n$  numbers. Suppose we want to determine, given a number  $x$ , the position of  $x$  in the array  $A$ , if any. One solution to the dictionary problem is to sort  $A$  (remembering every element’s original position) and then use binary search. The (implicit) binary search tree can be used almost directly as a decision tree. Each internal node in the *search tree* stores a key  $k$ ; the corresponding node in the *decision tree* stores the question ‘Is  $x < k$ ?’. Each leaf in the *search tree* stores some value  $A[i]$ ; the corresponding node in the *decision tree* asks ‘Is  $x = A[i]$ ?’ and has two leaf children, one labeled ‘ $i$ ’ and the other ‘none’.

We *define* the running time of a decision tree algorithm for a given input to be the number of queries in the path from the root to the leaf. For example, in the ‘Guess the animal’ tree above,  $T(\text{frog}) = 2$ . Thus, the worst-case running time of the algorithm is just the depth of the tree. This definition ignores other kinds of operations that the algorithm might perform that have nothing to do with the queries. (Even the most efficient binary search problem requires more than one machine instruction per comparison!)

<sup>2</sup>Complexity-theory snobs purists sometimes argue that ‘all algorithms’ is just a synonym for ‘all Turing machines’. This is utter nonsense; Turing machines are just another model of computation. Turing machines *might* be a reasonable abstraction of *physically realizable* computation—that’s the Church-Turing thesis—but it has a few problems. First, computation is an abstract mathematical process, not a physical process. Algorithms that use physically unrealistic components (like exact real numbers, or unbounded memory) are still mathematically well-defined and still provide useful intuition about real-world computation. Moreover, Turing machines don’t accurately reflect the complexity of physically realizable algorithms, because (for example) they can’t do arithmetic or access arbitrary memory locations in constant time. At best, they estimate algorithmic complexity up to polynomial factors (although even that is unknown).



But the number of decisions is certainly a *lower bound* on the actual running time, which is good enough to prove a lower bound on the complexity of a problem.

Both of the examples describe *binary* decision trees, where every query has only two answers. We may sometimes want to consider decision trees with higher degree. For example, we might use queries like ‘Is  $x$  greater than, equal to, or less than  $y$ ?’ or ‘Are these three points in clockwise order, colinear, or in counterclockwise order?’ A  $k$ -ary decision tree is one where every query has (at most)  $k$  different answers. **From now on, I will only consider  $k$ -ary decision trees where  $k$  is a constant.**

### 27.3 Information Theory

Most lower bounds for decision trees are based on the following simple observation: ***The answers to the queries must give you enough information to specify any possible output.*** If a problem has  $N$  different outputs, then obviously any decision tree must have at least  $N$  leaves. (It’s possible for several leaves to specify the same output.) Thus, if every query has at most  $k$  possible answers, then the depth of the decision tree must be at least  $\lceil \log_k N \rceil = \Omega(\log N)$ .

Let’s apply this to the dictionary problem for a set  $S$  of  $n$  numbers. Since there are  $n + 1$  possible outputs, any decision tree must have at least  $n + 1$  leaves, and thus any decision tree must have depth at least  $\lceil \log_k(n + 1) \rceil = \Omega(\log n)$ . So the complexity of the dictionary problem, in the decision-tree model of computation, is  $\Omega(\log n)$ . This matches the upper bound  $O(\log n)$  that comes from a perfectly-balanced binary search tree. That means that the standard binary search algorithm, which runs in  $O(\log n)$  time, is *optimal*—there is no faster algorithm in this model of computation.

### 27.4 But wait a second...

We can solve the membership problem in  $O(1)$  expected time using hashing. Isn’t this inconsistent with the  $\Omega(\log n)$  lower bound?

No, it isn’t. The reason is that hashing involves a query with more than a constant number of outcomes, specifically ‘What is the hash value of  $x$ ?’. In fact, if we don’t restrict the degree of the decision tree, we can get constant running time even without hashing, by using the obviously unreasonable query ‘For which index  $i$  (if any) is  $A[i] = x$ ?’. No, I am *not* cheating — remember that the decision tree model allows us to ask *any* question about the input!

This example illustrates a common theme in proving lower bounds: *choosing the right model of computation is absolutely crucial*. If you choose a model that is too powerful, the problem you’re studying may have a completely trivial algorithm. On the other hand, if you consider more restrictive models, the problem may not be solvable at all, in which case any lower bound will be meaningless! (In this class, we’ll just tell you the right model of computation to use.)

## 27.5 Sorting

Now let's consider the classical *sorting* problem — Given an array of  $n$  numbers, arrange them in increasing order. Unfortunately, decision trees don't have any way of describing moving data around, so we have to rephrase the question slightly:

Given a sequence  $\langle x_1, x_2, \dots, x_n \rangle$  of  $n$  distinct numbers, find the permutation  $\pi$  such that  $x_{\pi(1)} < x_{\pi(2)} < \dots < x_{\pi(n)}$ .

Now a  $k$ -ary decision-tree lower bound is immediate. Since there are  $n!$  possible permutations  $\pi$ , any decision tree for sorting must have at least  $n!$  leaves, and so must have depth  $\Omega(\log(n!))$ . To simplify the lower bound, we apply *Stirling's approximation*

$$n! = \left(\frac{n}{e}\right)^n \sqrt{2\pi n} \left(1 + \Theta\left(\frac{1}{n}\right)\right) > \left(\frac{n}{e}\right)^n.$$

This gives us the lower bound

$$\lceil \log_k(n!) \rceil > \left\lceil \log_k \left(\frac{n}{e}\right)^n \right\rceil = \lceil n \log_k n - n \log_k e \rceil = \Omega(n \log n).$$

This matches the  $O(n \log n)$  upper bound that we get from mergesort, heapsort, or quicksort, so those algorithms are optimal. The decision-tree complexity of sorting is  $\Theta(n \log n)$ .

Well... we're not quite done. In order to say that those algorithms are optimal, we have to demonstrate that they fit into our model of computation. A few minutes of thought will convince you that they can be described as a special type of decision tree called a *comparison tree*, where every query is of the form 'Is  $x_i$  bigger or smaller than  $x_j$ ?' These algorithms treat any two input sequences exactly the same way as long as the same comparisons produce exactly the same results. This is a feature of any comparison tree. In other words, *the actual input values don't matter, only their order*. Comparison trees describe almost all well-known sorting algorithms: bubble sort, selection sort, insertion sort, shell sort, quicksort, heapsort, mergesort, and so forth—but *not* radix sort or bucket sort.

## 27.6 Finding the Maximum and Adversaries

Finally let's consider the **maximum problem**: Given an array  $X$  of  $n$  numbers, find its largest entry. Unfortunately, there's no hope of proving a lower bound in this formulation, since there are an infinite number of possible answers, so let's rephrase it slightly.

Given a sequence  $\langle x_1, x_2, \dots, x_n \rangle$  of  $n$  distinct numbers, find the index  $m$  such that  $x_m$  is the largest element in the sequence.

We can get an upper bound of  $n - 1$  comparisons in several different ways. The easiest is probably to start at one end of the sequence and do a linear scan, maintaining a current maximum. Intuitively, this seems like the best we can do, but the information-theoretic bound is only  $\lceil \log_2 n \rceil$ . And in fact, this bound is tight! We can locate the maximum element by asking only  $\lceil \log_2 n \rceil$  'unreasonable' questions like "Is the index of the maximum element odd?" No, this is *not* cheating—the decision tree model allows *arbitrary* questions.

To prove a non-trivial lower bound for this problem, we must do two things. First, we need to consider a more reasonable model of computation, by restricting the kinds of questions the algorithm is allowed to ask. We will consider the **comparison tree model**, where every query must have the form "Is

$x_i > x_j$ ?". Since most algorithms<sup>3</sup> for finding the maximum rely on comparisons to make control-flow decisions, this does not seem like an unreasonable restriction.

Second, we will use something called an **adversary argument**. The idea is that an all-powerful malicious adversary *pretends* to choose an input for the algorithm. When the algorithm asks a question about the input, the adversary answers in whatever way will make the algorithm do the most work. If the algorithm does not ask enough queries before terminating, then there will be several different inputs, each consistent with the adversary's answers, that should result in different outputs. In this case, whatever the algorithm outputs, the adversary can 'reveal' an input that is consistent with its answers, but contradicts the algorithm's output, and then claim that that was the input that he was using all along.

For the maximum problem, the adversary originally pretends that  $x_i = i$  for all  $i$ , and answers all comparison queries accordingly. Whenever the adversary reveals that  $x_i < x_j$ , he *marks*  $x_i$  as an item that the algorithm knows (or should know) is not the maximum element. At most one element  $x_i$  is marked after each comparison. Note that  $x_n$  is never marked. If the algorithm does less than  $n - 1$  comparisons before it terminates, the adversary must have at least one other unmarked element  $x_k \neq x_n$ . In this case, the adversary can change the value of  $x_k$  from  $k$  to  $n + 1$ , making  $x_k$  the largest element, without being inconsistent with any of the comparisons that the algorithm has performed. In other words, the algorithm cannot tell that the adversary has cheated. However,  $x_n$  is the maximum element in the original input, and  $x_k$  is the largest element in the modified input, so the algorithm cannot possibly give the correct answer for both cases. Thus, in order to be correct, any algorithm must perform at least  $n - 1$  comparisons.

The adversary argument we described has two very important properties. First, no algorithm can distinguish between a malicious adversary and an honest user who actually chooses an input in advance and answers all queries truthfully. But much more importantly, **the adversary makes absolutely no assumptions about the order in which the algorithm performs comparisons**. The adversary forces *any* comparison-based algorithm<sup>4</sup> to either perform  $n - 1$  comparisons, or to give the wrong answer for at least one input sequence.

## Exercises

0. Simon bar Kokhba thinks of an integer between 1 and 1,000,000 (or so he claims). You are trying to determine his number by asking as few yes/no questions as possible. How many yes/no questions are required to determine Simon's number in the worst case? Give both an upper bound (supported by an algorithm) and a lower bound.
  
  
  
  
  
  
  
  
1. Consider the following *multi-dictionary* problem. Let  $A[1..n]$  be a fixed array of distinct integers. Given an array  $X[1..k]$ , we want to find the position (if any) of each integer  $X[i]$  in the array  $A$ . In other words, we want to compute an array  $I[1..k]$  where for each  $i$ , either  $I[i] = 0$  (so zero means 'none') or  $A[I[i]] = X[i]$ . Determine the *exact* complexity of this problem, as a function of  $n$  and  $k$ , in the binary decision tree model.

---

<sup>3</sup>but not all—see Exercise ??

<sup>4</sup>In fact, the  $n - 1$  lower bound for finding the maximum holds in a more powerful model called *algebraic* decision trees, which are binary trees where every query is a comparison between two polynomial functions of the input values, such as 'Is  $x_1^2 - 3x_2x_3 + x_4^{17}$  bigger or smaller than  $5 + x_1x_3^5x_5^2 - 2x_7^{42}$ ?'

2. We say that an array  $A[1..n]$  is *k-sorted* if it can be divided into  $k$  blocks, each of size  $n/k$ , such that the elements in each block are larger than the elements in earlier blocks, and smaller than elements in later blocks. The elements within each block need not be sorted.

For example, the following array is 4-sorted:

1	2	4	3	7	6	8	5	10	11	9	12	15	13	16	14
---	---	---	---	---	---	---	---	----	----	---	----	----	----	----	----

- (a) Describe an algorithm that *k*-sorts an arbitrary array in  $O(n \log k)$  time.
- (b) Prove that any comparison-based *k*-sorting algorithm requires  $\Omega(n \log k)$  comparisons in the worst case.
- (c) Describe an algorithm that completely sorts an already *k*-sorted array in  $O(n \log(n/k))$  time.
- (d) Prove that any comparison-based algorithm to completely sort a *k*-sorted array requires  $\Omega(n \log(n/k))$  comparisons in the worst case.

In all cases, you can assume that  $n/k$  is an integer.

3. Recall the nuts-and-bolts problem from the lecture on randomized algorithms. We are given  $n$  bolts and  $n$  nuts of different sizes, where each bolt exactly matches one nut. Our goal is to find the matching nut for each bolt. The nuts and bolts are too similar to compare directly; however, we can test whether any nut is too big, too small, or the same size as any bolt.
- (a) Prove that in the worst case,  $\Omega(n \log n)$  nut-bolt tests are required to correctly match up the nuts and bolts.
  - (b) Now suppose we would be happy to find *most* of the matching pairs. Prove that in the worst case,  $\Omega(n \log n)$  nut-bolt tests are required even to find  $n/2$  arbitrary matching nut-bolt pairs.
  - \*(c) Prove that in the worst case,  $\Omega(n + k \log n)$  nut-bolt tests are required to find  $k$  arbitrary matching pairs. [*Hint: Use an adversary argument for the  $\Omega(n)$  term.*]
  - \*(d) Describe a randomized algorithm that finds  $k$  matching nut-bolt pairs in  $O(n + k \log n)$  expected time.
- \*4. Suppose you want to determine the largest number in an  $n$ -element set  $X = \{x_1, x_2, \dots, x_n\}$ , where each element  $x_i$  is an integer between 1 and  $2^m - 1$ . Describe an algorithm that solves this problem in  $O(n + m)$  steps, where at each step, your algorithm compares one of the elements  $x_i$  with a *constant*. In particular, your algorithm must never actually compare two elements of  $X$ ! [*Hint: Construct and maintain a nested set of ‘pinning intervals’ for the numbers that you have not yet removed from consideration, where each interval but the largest is either the upper half or lower half of the next larger block.*]

*An adversary means opposition and competition,  
but not having an adversary means grief and loneliness.*

— Zhuangzi (Chuang-tsu) c. 300 BC

*It is possible that the operator could be hit by an asteroid and your \$20 could fall off his cardboard box and land on the ground, and while you were picking it up, \$5 could blow into your hand. You therefore could win \$5 by a simple twist of fate.*

— Penn Jillette, explaining how to win at Three-Card Monte (1999)

## 28 Adversary Arguments

### 28.1 Three-Card Monte

Until Times Square was sanitized into TimesSquareLand™, you could often find dealers stealing tourists' money using a game called 'Three Card Monte' or 'Spot the Lady'. The dealer has three cards, say the Queen of Hearts, the two of spades, and three of clubs. The dealer shuffles the cards face down on a table (usually slowly enough that you can follow the Queen), and then asks the tourist to bet on which card is the Queen. In principle, the tourist's odds of winning are at least one in three.

In practice, however, the tourist *never* wins, because the dealer cheats. The dealer actually holds at least *four* cards; before he even starts shuffling the cards, the dealer palms the queen or sticks it up his sleeve. No matter what card the tourist bets on, the dealer turns over a black card (which might be the two of clubs, but most tourists won't notice that wasn't one of the original cards). If the tourist gives up, the dealer slides the queen under one of the cards and turns it over, showing the tourist 'where the queen was all along'. If the dealer is really good, the tourist won't see the dealer changing the cards and will think maybe the queen *was* there all along and he just wasn't smart enough to figure that out. As long as the dealer doesn't reveal all the black cards at once, the tourist has no way to prove that the dealer cheated!<sup>1</sup>

### 28.2 *n*-Card Monte

Now let's consider a similar game, but with an algorithm acting as the tourist and with bits instead of cards. Suppose we have an array of  $n$  bits and we want to determine if any of them is a 1. Obviously we can figure this out by just looking at every bit, but can we do better? Is there maybe some complicated tricky algorithm to answer the question "Any ones?" without looking at every bit? Well, of course not, but how do we prove it?

The simplest proof technique is called an *adversary* argument. The idea is that an all-powerful malicious adversary (the dealer) *pretends* to choose an input for the algorithm (the tourist). When the algorithm wants looks at a bit (a card), the adversary sets that bit to whatever value will make the algorithm do the most work. If the algorithm does not look at enough bits before terminating, then there will be several different inputs, each consistent with the bits already seen, the should result in different outputs. Whatever the algorithm outputs, the adversary can 'reveal' an input that it has all the examined bits but contradicts the algorithm's output, and then claim that that was the input that he was using all along. Since the only information the algorithm has is the set of bits it examined, the algorithm

---

<sup>1</sup>Even if the dealer is a sloppy magician, he'll cheat anyway. The dealer is almost always surrounded by shills; these are the 'tourists' who look like they're actually winning, who turn over cards when the dealer 'isn't looking', who casually mention how easy the game is to win, and so on. The shills protect the dealer from any angry tourists who notice the dealer cheating, and shake down any tourists who refuse to pay after making a bet. You **cannot** win this game, **ever**.

cannot distinguish between a malicious adversary and an honest user who actually chooses an input in advance and answers all queries truthfully.

For the  $n$ -card monte problem, the adversary originally pretends that the input array is all zeros—whenever the algorithm looks at a bit, it sees a 0. Now suppose the algorithm stops before looking at all three bits. If the algorithm says ‘No, there’s no 1,’ the adversary changes one of the unexamined bits to a 1 and shows the algorithm that it’s wrong. If the algorithm says ‘Yes, there’s a 1,’ the adversary reveals the array of zeros and again proves the algorithm wrong. Either way, the algorithm cannot tell that the adversary has cheated.

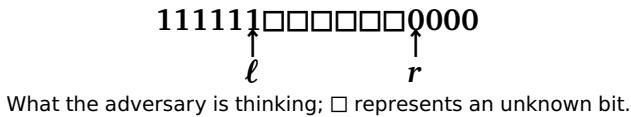
One absolutely crucial feature of this argument is that *the adversary makes absolutely no assumptions about the algorithm*. The adversary strategy can’t depend on some predetermined order of examining bits, and it doesn’t care about anything the algorithm might or might not do when it’s not looking at bits. Any algorithm that doesn’t examine every bit falls victim to the adversary.

### 28.3 Finding Patterns in Bit Strings

Let’s make the problem a little more complicated. Suppose we’re given an array of  $n$  bits and we want to know if it contains the substring 01, a zero followed immediately by a one. Can we answer this question without looking at every bit?

It turns out that if  $n$  is odd, we *don’t* have to look at all the bits. First we look at the bits in every even position:  $B[2], B[4], \dots, B[n - 1]$ . If we see  $B[i] = 0$  and  $B[j] = 1$  for any  $i < j$ , then we know the pattern 01 is in there somewhere—starting at the last 0 before  $B[j]$ —so we can stop without looking at any more bits. If we see only 1s followed by 0s, we don’t have to look at the bit between the last 0 and the first 1. If every even bit is a 0, we don’t have to look at  $B[1]$ , and if every even bit is a 1, we don’t have to look at  $B[n]$ . In the worst case, our algorithm looks at only  $n - 1$  of the  $n$  bits.

But what if  $n$  is even? In that case, we can use the following adversary strategy to show that any algorithm *does* have to look at every bit. The adversary will attempt to produce an ‘input’ string  $B$  *without* the substring 01; all such strings have the form 11...100...0. The adversary maintains two indices  $\ell$  and  $r$  and pretends that the prefix  $B[1.. \ell]$  contains only 1s and the suffix  $B[r..n]$  contains only 0s. Initially  $\ell = 0$  and  $r = n + 1$ .



What the adversary is thinking; □ represents an unknown bit.

The adversary maintains the invariant that  $r - \ell$ , the length of the undecided portion of the ‘input’ string, is even. When the algorithm looks at a bit between  $\ell$  and  $r$ , the adversary chooses whichever value preserves the parity of the intermediate chunk of the array, and then moves either  $\ell$  or  $r$ . Specifically, here’s what the adversary does when the algorithm examines bit  $B[i]$ . (Note that I’m specifying the adversary strategy as an algorithm!)

```

HIDE01( $i$ ):
  if  $i \leq \ell$ 
     $B[i] \leftarrow 1$ 
  else if  $i \geq r$ 
     $B[i] \leftarrow 0$ 
  else if  $i - \ell$  is even
     $B[i] \leftarrow 0$ 
     $r \leftarrow i$ 
  else
     $B[i] \leftarrow 1$ 
     $\ell \leftarrow i$ 

```

It's fairly easy to prove that this strategy forces the algorithm to examine every bit. If the algorithm doesn't look at every bit to the right of  $r$ , the adversary could replace some unexamined bit with a 1. Similarly, if the algorithm doesn't look at every bit to the left of  $\ell$ , the adversary could replace some unexamined bit with a zero. Finally, if there are any unexamined bits between  $\ell$  and  $r$ , there must be at least two such bits (since  $r - \ell$  is always even) and the adversary can put a 01 in the gap.

In general, we say that a bit pattern is *evasive* if we have to look at every bit to decide if a string of  $n$  bits contains the pattern. So the pattern 1 is evasive for all  $n$ , and the pattern 01 is evasive if and only if  $n$  is even. It turns out that the *only* patterns that are evasive for *all* values of  $n$  are the one-bit patterns 0 and 1.

## 28.4 Evasive Graph Properties

Another class of problems for which adversary arguments give good lower bounds is graph problems where the graph is represented by an adjacency matrix, rather than an adjacency list. Recall that the adjacency matrix of an undirected  $n$ -vertex graph  $G = (V, E)$  is an  $n \times n$  matrix  $A$ , where  $A[i, j] = \llbracket (i, j) \in E \rrbracket$ . We are interested in deciding whether an undirected graph has or does not have a certain *property*. For example, is the input graph connected? Acyclic? Planar? Complete? A tree? We call a graph property *evasive* if we have to look at all  $\binom{n}{2}$  entries in the adjacency matrix to decide whether a graph has that property.

An obvious example of an evasive graph property is *emptiness*: Does the graph have any edges at all? We can show that emptiness is evasive using the following simple adversary strategy. The adversary maintains two graphs  $E$  and  $G$ .  $E$  is just the empty graph with  $n$  vertices. Initially  $G$  is the complete graph on  $n$  vertices. Whenever the algorithm asks about an edge, the adversary removes that edge from  $G$  (unless it's already gone) and answers 'no'. If the algorithm terminates without examining every edge, then  $G$  is not empty. Since both  $G$  and  $E$  are consistent with all the adversary's answers, the algorithm must give the wrong answer for one of the two graphs.

## 28.5 Connectedness Is Evasive

Now let me give a more complicated example, *connectedness*. Once again, the adversary maintains two graphs,  $Y$  and  $M$  ('yes' and 'maybe').  $Y$  contains all the edges that the algorithm knows are definitely in the input graph.  $M$  contains all the edges that the algorithm thinks *might* be in the input graph, or in other words, all the edges of  $Y$  plus all the unexamined edges. Initially,  $Y$  is empty and  $M$  is complete.

Here's the strategy that adversary follows when the adversary asks whether the input graph contains the edge  $e$ . I'll assume that whenever an algorithm examines an edge, it's in  $M$  but not in  $Y$ ; in other words, algorithms never ask about the same edge more than once.

```

HIDECONNECTEDNESS( $e$ ):
  if  $M \setminus \{e\}$  is connected
    remove  $(i, j)$  from  $M$ 
    return 0
  else
    add  $e$  to  $Y$ 
    return 1

```

Notice that the graphs  $Y$  and  $M$  are both consistent with the adversary's answers at all times. The adversary strategy maintains a few other simple invariants.

- **$Y$  is a subgraph of  $M$ .** This is obvious.
- **$M$  is connected.** This is also obvious.
- **If  $M$  has a cycle, none of its edges are in  $Y$ .** If  $M$  has a cycle, then deleting any edge in that cycle leaves  $M$  connected.
- **$Y$  is acyclic.** This follows directly from the previous invariant.
- **If  $Y \neq M$ , then  $Y$  is disconnected.** The only connected acyclic graph is a tree. Suppose  $Y$  is a tree and some edge  $e$  is in  $M$  but not in  $Y$ . Then there is a cycle in  $M$  that contains  $e$ , all of whose other edges are in  $Y$ . This violated our third invariant.

We can also think about the adversary strategy in terms of minimum spanning trees. Recall the anti-Kruskal algorithm for computing the *maximum* spanning tree of a graph: Consider the edges one at a time in increasing order of length. If removing an edge would disconnect the graph, declare it part of the spanning tree (by adding it to  $Y$ ); otherwise, throw it away (by removing it from  $M$ ). If the algorithm examines all  $\binom{n}{2}$  possible edges, then  $Y$  and  $M$  are both equal to the maximum spanning tree of the complete  $n$ -vertex graph, where the weight of an edge is the time when the algorithm asked about it.

Now, if an algorithm terminates before examining all  $\binom{n}{2}$  edges, then there is at least one edge in  $M$  that is not in  $Y$ . Since the algorithm cannot distinguish between  $M$  and  $Y$ , even though  $M$  is connected and  $Y$  is not, the algorithm cannot possibly give the correct output for both graphs. Thus, in order to be correct, any algorithm must examine every edge—*Connectedness is evasive!*

## 28.6 An Evasive Conjecture

A graph property is *nontrivial* if there is at least one graph with the property and at least one graph without the property. (The only trivial properties are ‘Yes’ and ‘No’.) A graph property is *monotone* if it is closed under taking subgraphs — if  $G$  has the property, then any subgraph of  $G$  has the property. For example, emptiness, planarity, acyclicity, and *non-connectedness* are monotone. The properties of being a tree and of having a vertex of degree 3 are not monotone.

**Conjecture 1 (Aanderraa, Karp, and Rosenberg).** *Every nontrivial monotone property of  $n$ -vertex graphs is evasive.*

The Aanderraa-Karp-Rosenberg conjecture has been proven when  $n = p^e$  for some prime  $p$  and positive integer exponent  $e$ —the proof uses some interesting results from algebraic topology<sup>2</sup>—but it is

<sup>2</sup>Let  $\Delta$  be a contractible simplicial complex whose automorphism group  $\text{Aut}(\Delta)$  is vertex-transitive, and let  $\Gamma$  be a vertex-transitive subgroup of  $\text{Aut}(\Delta)$ . If there are normal subgroups  $\Gamma_1 \triangleleft \Gamma_2 \triangleleft \Gamma$  such that  $|\Gamma_1| = p^\alpha$  for some prime  $p$  and integer  $\alpha$ ,  $|\Gamma/\Gamma_2| = q^\beta$  for some prime  $q$  and integer  $\beta$ , and  $\Gamma_2/\Gamma_1$  is cyclic, then  $\Delta$  is a simplex.

No, this will not be on the final exam.

still open for other values of  $n$ .

There are non-trivial non-evasive graph properties, but all known examples are non-monotone. One such property—‘scorpionhood’—is described in an exercise at the end of this lecture note.

## 28.7 Finding the Minimum and Maximum

Last time, we saw an adversary argument that finding the largest element of an unsorted set of  $n$  numbers requires at least  $n - 1$  comparisons. Let’s consider the complexity of finding the largest *and* smallest elements. More formally:

Given a sequence  $X = \langle x_1, x_2, \dots, x_n \rangle$  of  $n$  distinct numbers, find indices  $i$  and  $j$  such that  $x_i = \min X$  and  $x_j = \max X$ .

How many comparisons do we need to solve this problem? An upper bound of  $2n - 3$  is easy: find the minimum in  $n - 1$  comparisons, and then find the maximum of everything else in  $n - 2$  comparisons. Similarly, a lower bound of  $n - 1$  is easy, since any algorithm that finds the min and the max certainly finds the max.

We can improve both the upper and the lower bound to  $\lceil 3n/2 \rceil - 2$ . The upper bound is established by the following algorithm. Compare all  $\lfloor n/2 \rfloor$  consecutive pairs of elements  $x_{2i-1}$  and  $x_{2i}$ , and put the smaller element into a set  $S$  and the larger element into a set  $L$ . If  $n$  is odd, put  $x_n$  into both  $L$  and  $S$ . Then find the smallest element of  $S$  and the largest element of  $L$ . The total number of comparisons is at most

$$\underbrace{\left\lfloor \frac{n}{2} \right\rfloor}_{\text{build } S \text{ and } L} + \underbrace{\left\lceil \frac{n}{2} \right\rceil - 1}_{\text{compute } \min S} + \underbrace{\left\lceil \frac{n}{2} \right\rceil - 1}_{\text{compute } \max L} = \left\lceil \frac{3n}{2} \right\rceil - 2.$$

For the lower bound, we use an adversary argument. The adversary marks each element + if it *might* be the maximum element, and – if it *might* be the minimum element. Initially, the adversary puts both marks + and – on every element. If the algorithm compares two double-marked elements, then the adversary declares one smaller, removes the + mark from the smaller element, and removes the – mark from the larger one. In every other case, the adversary can answer so that at most one mark needs to be removed. For example, if the algorithm compares a double marked element against one labeled –, the adversary says the one labeled – is smaller and removes the – mark from the other. If the algorithm compares to +'s, the adversary must unmark one of the two.

Initially, there are  $2n$  marks. At the end, in order to be correct, exactly one item must be marked + and exactly one other must be marked –, since the adversary can make any + the maximum and any – the minimum. Thus, the algorithm must force the adversary to remove  $2n - 2$  marks. At most  $\lfloor n/2 \rfloor$  comparisons remove two marks; every other comparison removes at most one mark. Thus, the adversary strategy forces any algorithm to perform at least  $2n - 2 - \lfloor n/2 \rfloor = \lceil 3n/2 \rceil - 2$  comparisons.

## 28.8 Finding the Median

Finally, let’s consider the *median* problem: Given an unsorted array  $X$  of  $n$  numbers, find its  $n/2$ th largest entry. (I’ll assume that  $n$  is even to eliminate pesky floors and ceilings.) More formally:

Given a sequence  $\langle x_1, x_2, \dots, x_n \rangle$  of  $n$  distinct numbers, find the index  $m$  such that  $x_m$  is the  $n/2$ th largest element in the sequence.

To prove a lower bound for this problem, we can use a combination of information theory and two adversary arguments. We use one adversary argument to prove the following simple lemma:

**Lemma 1.** Any comparison tree that correctly finds the median element also identifies the elements smaller than the median and larger than the median.

**Proof:** Suppose we reach a leaf of a decision tree that chooses the median element  $x_m$ , and there is still some element  $x_i$  that isn't known to be larger or smaller than  $x_m$ . In other words, we cannot decide based on the comparisons that we've already performed whether  $x_i < x_m$  or  $x_i > x_m$ . Then in particular no element is known to lie between  $x_i$  and  $x_m$ . This means that there must be an input that is consistent with the comparisons we've performed, in which  $x_i$  and  $x_m$  are adjacent in sorted order. But then we can swap  $x_i$  and  $x_m$ , without changing the result of any comparison, and obtain a different consistent input in which  $x_i$  is the median, not  $x_m$ . Our decision tree gives the wrong answer for this 'swapped' input.  $\square$

This lemma lets us rephrase the median-finding problem yet again.

Given a sequence  $X = \langle x_1, x_2, \dots, x_n \rangle$  of  $n$  distinct numbers, find the indices of its  $n/2 - 1$  largest elements  $L$  and its  $n/2$ th largest element  $x_m$ .

Now suppose a 'little birdie' tells us the set  $L$  of elements larger than the median. This information fixes the outcomes of certain comparisons—any item in  $L$  is bigger than any element not in  $L$ —so we can 'prune' those comparisons from the comparison tree. The pruned tree finds the largest element of  $X \setminus L$  (the median of  $X$ ), and thus must have depth at least  $n/2 - 1$ . In fact, the adversary argument in the last lecture implies that every leaf in the pruned tree must have depth at least  $n/2 - 1$ , so the pruned tree has at least  $2^{n/2-1}$  leaves.

There are  $\binom{n}{n/2-1} \approx 2^n / \sqrt{n/2}$  possible choices for the set  $L$ . Every leaf in the original comparison tree is also a leaf in exactly one of the  $\binom{n}{n/2-1}$  pruned trees, so the original comparison tree must have at least  $\binom{n}{n/2-1} 2^{n/2-1} \approx 2^{3n/2} / \sqrt{n/2}$  leaves. Thus, any comparison tree that finds the median must have depth at least

$$\left\lceil \frac{n}{2} - 1 + \lg \binom{n}{n/2-1} \right\rceil = \frac{3n}{2} - O(\log n).$$

A more complicated adversary argument (also involving pruning the comparison tree with little birdies) improves this lower bound to  $2n - o(n)$ .

A similar argument implies that at least  $n - k + \lceil \lg \binom{n}{k-1} \rceil = \Omega((n - k) + k \log(n/k))$  comparisons are required to find the  $k$ th largest element in an  $n$ -element set. This bound is tight up to constant factors for all  $k \leq n/2$ ; there is an algorithm that uses at most  $O(n + k \log(n/k))$  comparisons. Moreover, this lower bound is *exactly* tight when  $k = 1$  or  $k = 2$ . In fact, these are the *only* values of  $k \leq n/2$  for which the exact complexity of the selection problem is known. Even the case  $k = 3$  is still open!

## Exercises

1. (a) Let  $X$  be a set containing an odd number of  $n$ -bit strings. Prove that any algorithm that decides whether a given  $n$ -bit string is an element of  $X$  *must* examine every bit of the input string in the worst case.
- (b) Give a one-line proof that the bit pattern 01 is evasive for all even  $n$ .
- (c) Prove that the bit pattern 11 is evasive if and only if  $n \bmod 3 = 1$ .
- \*(d) Prove that the bit pattern 111 is evasive if and only if  $n \bmod 4 = 0$  or 3.

2. Suppose we are given the adjacency matrix of a *directed* graph  $G$  with  $n$  vertices. Describe an algorithm that determines whether  $G$  has a *sink* by probing only  $O(n)$  bits in the input matrix. A sink is a vertex that has an incoming edge from every other vertex, but no outgoing edges.
  
- \*3. A *scorpion* is an undirected graph with three special vertices: the *sting*, the *tail*, and the *body*. The sting is connected only to the tail; the tail is connected only to the sting and the body; and the body is connected to every vertex except the sting. The rest of the vertices (the head, eyes, legs, antennae, teeth, gills, flippers, wheels, etc.) can be connected arbitrarily. Describe an algorithm that determines whether a given  $n$ -vertex graph is a scorpion by probing only  $O(n)$  entries in the adjacency matrix.
  
4. Prove using an adversary argument that acyclicity is an evasive graph property. [Hint: Kruskal.]
  
5. Prove that finding the second largest element in an  $n$ -element array requires *exactly*  $n - 2 + \lceil \lg n \rceil$  comparisons in the worst case. Prove the upper bound by describing and analyzing an algorithm; prove the lower bound using an adversary argument.
  
6. Let  $T$  be a perfect ternary tree where every leaf has depth  $\ell$ . Suppose each of the  $3^\ell$  leaves of  $T$  is labeled with a bit, either 0 or 1, and each internal node is labeled with a bit that agrees with the *majority* of its children.
  - (a) Prove that any deterministic algorithm that determines the label of the root must examine all  $3^\ell$  leaf bits in the worst case.
  - (b) Describe and analyze a *randomized* algorithm that determines the root label, such that the expected number of leaves examined is  $o(3^\ell)$ . (You may want to review the notes on randomized algorithms.)
  
- \*7. UIUC has just finished constructing the new Reingold Building, the tallest dormitory on campus. In order to determine how much insurance to buy, the university administration needs to determine the highest safe floor in the building. A floor is considered *safe* if a ~~drunk student~~ *an egg* can fall from a window on that floor and land without breaking; if the egg breaks, the floor is considered *unsafe*. Any floor that is higher than an unsafe floor is also considered unsafe. The only way to determine whether a floor is safe is to drop an egg from a window on that floor.

You would like to find the lowest unsafe floor  $L$  by performing as few tests as possible; unfortunately, you have only a very limited supply of eggs.

- (a) Prove that if you have only one egg, you can find the lowest unsafe floor with  $L$  tests. [Hint: Yes, this is trivial.]
- (b) Prove that if you have only one egg, you must perform at least  $L$  tests in the worst case. In other words, prove that your algorithm from part (a) is optimal. [Hint: Use an adversary argument.]
- (c) Describe an algorithm to find the lowest unsafe floor using *two* eggs and only  $O(\sqrt{L})$  tests. [Hint: Ideally, each egg should be dropped the same number of times. How many floors can you test with  $n$  drops?]

- (d) Prove that if you start with two eggs, you must perform at least  $\Omega(\sqrt{L})$  tests in the worst case. In other words, prove that your algorithm from part (c) is optimal.
- \*(e) Describe an algorithm to find the lowest unsafe floor using  $k$  eggs, using as few tests as possible, and prove your algorithm is optimal for all values of  $k$ .

*The wonderful thing about standards is that there are so many of them to choose from.*

— Real Admiral Grace Murray Hopper

*If a problem has no solution, it may not be a problem, but a fact — not to be solved, but to be coped with over time.*

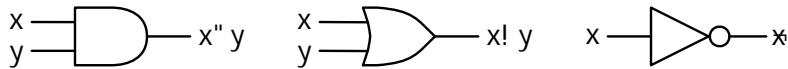
— Shimon Peres

## 29 NP-Hard Problems

### 29.1 ‘Efficient’ Problems

A generally-accepted minimum requirement for an algorithm to be considered ‘efficient’ is that its running time is polynomial:  $O(n^c)$  for some constant  $c$ , where  $n$  is the size of the input.<sup>1</sup> Researchers recognized early on that not all problems can be solved this quickly, but we had a hard time figuring out exactly which ones could and which ones couldn’t. There are several so-called *NP-hard* problems, which most people believe *cannot* be solved in polynomial time, even though nobody can prove a super-polynomial lower bound.

*Circuit satisfiability* is a good example of a problem that we don’t know how to solve in polynomial time. In this problem, the input is a *boolean circuit*: a collection of AND, OR, and NOT gates connected by wires. We will assume that there are no loops in the circuit (so no delay lines or flip-flops). The input to the *circuit* is a set of  $m$  boolean (TRUE/FALSE) values  $x_1, \dots, x_m$ . The output is a single boolean value. Given specific input values, we can calculate the output of the circuit in polynomial (actually, *linear*) time using depth-first-search, since we can compute the output of a  $k$ -input gate in  $O(k)$  time.



An AND gate, an OR gate, and a NOT gate.

A boolean circuit. Inputs enter from the left, and the output leaves to the right.

The circuit satisfiability problem asks, given a circuit, whether there is an input that makes the circuit output TRUE, or conversely, whether the circuit *always* outputs FALSE. Nobody knows how to solve

<sup>1</sup>This notion of efficiency was independently formalized by Alan Cobham (*The intrinsic computational difficulty of functions. Logic, Methodology, and Philosophy of Science (Proc. Int. Congress)*, 24–30, 1965), Jack Edmonds (*Paths, trees, and flowers. Canadian Journal of Mathematics* 17:449–467, 1965), and Michael Rabin (*Mathematical theory of automata. Proceedings of the 19th ACM Symposium in Applied Mathematics*, 153–175, 1966), although similar notions were considered more than a decade earlier by Kurt Gödel and John von Neumann.

this problem faster than just trying all  $2^m$  possible inputs to the circuit, but this requires exponential time. On the other hand, nobody has ever proved that this is the best we can do; maybe there's a clever algorithm that nobody has discovered yet!

## 29.2 P, NP, and co-NP

A *decision problem* is a problem whose output is a single boolean value: YES or NO.<sup>2</sup> Let me define three classes of decision problems:

- **P** is the set of decision problems that can be solved in polynomial time.<sup>3</sup> Intuitively, P is the set of problems that can be solved quickly.
- **NP** is the set of decision problems with the following property: If the answer is YES, then there is a *proof* of this fact that can be checked in polynomial time. Intuitively, NP is the set of decision problems where we can verify a YES answer quickly if we have the solution in front of us.
- **co-NP** is the opposite of NP. If the answer to a problem in co-NP is No, then there is a proof of this fact that can be checked in polynomial time.

For example, the circuit satisfiability problem is in NP. If the answer is YES, then any set of  $m$  input values that produces TRUE output is a proof of this fact; we can check the proof by evaluating the circuit in polynomial time. It is widely believed that circuit satisfiability is *not* in P or in co-NP, but nobody actually knows.

Every decision problem in P is also in NP. If a problem is in P, we can verify YES answers in polynomial time recomputing the answer from scratch! Similarly, any problem in P is also in co-NP.

Perhaps the single most important open questions in theoretical computer science, if not all of mathematics, is whether the complexity classes P and NP are actually different. Intuitively, it seems obvious to most people that  $P \neq NP$ ; the homeworks and exams in this class and others have (I hope) convinced you that problems can be incredibly hard to solve, even when the solutions are obvious in retrospect. The prominent physicist Richard Feynman, one of the early pioneers of quantum computing, apparently had trouble accepting that  $P \neq NP$  was actually an open question. But nobody knows how to prove it. The Clay Mathematics Institute lists P versus NP as the first of its seven Millennium Prize Problems, offering a \$1,000,000 reward for its solution.

A more subtle but still open question is whether the complexity classes NP and co-NP are different. Even if we can verify every YES answer quickly, there's no reason to believe we can also verify No answers quickly. For example, as far as we know, there is no short proof that a boolean circuit is *not* satisfiable. It is generally believed that  $NP \neq co\text{-}NP$ , but nobody knows how to prove it.

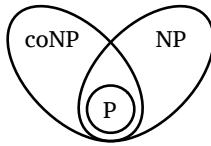
## 29.3 NP-hard, NP-easy, and NP-complete

A problem  $\Pi$  is **NP-hard** if a polynomial-time algorithm for  $\Pi$  would imply a polynomial-time algorithm for *every problem in NP*. In other words:

$$\boxed{\Pi \text{ is NP-hard} \iff \text{If } \Pi \text{ can be solved in polynomial time, then } P=NP}$$

<sup>2</sup>Technically, I should be talking about *languages*, which are just sets of bit strings. The language associated with a decision problem is the set of bit strings for which the answer is YES. For example, for the problem is ‘Is the input graph connected?’, the corresponding language is the set of connected graphs, where each graph is represented as a bit string (for example, its adjacency matrix).

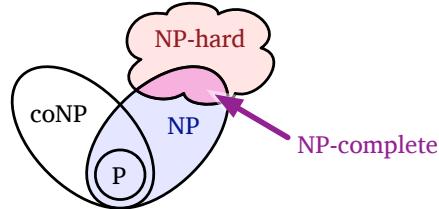
<sup>3</sup>More formally, P is the set of languages that can be recognized in polynomial time by a single-tape Turing machine. If you want to learn more about Turing machines, take CS 579.

What we *think* the world looks like.

Intuitively, if we could solve one particular NP-hard problem quickly, then we could quickly solve *any* problem whose solution is easy to understand, using the solution to that one special problem as a subroutine. NP-hard problems are at least as hard as any problem in NP.

Calling a problem is NP-hard is like saying ‘If I own a dog, then it can speak fluent English.’ You probably don’t know whether or not I own a dog, but you’re probably pretty sure that I don’t own a *talking* dog. Nobody has a mathematical *proof* that dogs can’t speak English—the fact that no one has ever heard a dog speak English is evidence, as are the hundreds of examinations of dogs that lacked the proper mouth shape and brainpower, but mere evidence is not a mathematical proof. Nevertheless, no sane person would believe me if I said I owned a dog that spoke fluent English. So the statement ‘If I own a dog, then it can speak fluent English’ has a natural corollary: No one in their right mind should believe that I own a dog! Likewise, if a problem is NP-hard, no one in their right mind should believe it can be solved in polynomial time.

Finally, a problem is *NP-complete* if it is both NP-hard and an element of NP (or ‘NP-easy’). NP-complete problems are the hardest problems in NP. If anyone finds a polynomial-time algorithm for even one NP-complete problem, then that would imply a polynomial-time algorithm for *every* NP-complete problem. Literally *thousands* of problems have been shown to be NP-complete, so a polynomial-time algorithm for one (and therefore all) of them seems incredibly unlikely.

More of what we *think* the world looks like.

It is not immediately clear that *any* decision problems are NP-hard or NP-complete. NP-hardness is already a lot to demand of a problem; insisting that the problem also have a nondeterministic polynomial-time algorithm seems almost completely unreasonable. The following remarkable theorem was first published by Steve Cook in 1971 and independently by Leonid Levin in 1973.<sup>4</sup> I won’t even sketch the proof, since I’ve been (deliberately) vague about the definitions.

**The Cook-Levin Theorem.** *Circuit satisfiability is NP-complete.*

#### \*29.4 Formal Definition (HC SVNT DRACONES)

More formally, a problem  $\Pi$  is defined to be NP-hard if and only if, for every problem  $\Pi'$  in NP, there is a polynomial-time **Turing reduction** from  $\Pi'$  to  $\Pi$ —a Turing reduction just means a reduction that can

<sup>4</sup>Levin first reported his results at seminars in Moscow in 1971; news of Cook’s result did not reach the Soviet Union until at least 1973, after Levin’s announcement of his results had been published. In accordance with Stigler’s Law, this result is often called ‘Cook’s Theorem’. Cook won the Turing award for his proof. Levin was denied his PhD at Moscow University for political reasons; he emigrated to the US in 1978 and earned a PhD at MIT a year later.

be executed on a Turing machine. Polynomial-time Turing reductions are also called *oracle reductions* or *Cook reductions*. It is elementary, but *extremely* tedious, to prove that any algorithm that can be executed on a random-access machine<sup>5</sup> in time  $T(n)$  can be simulated on a single-tape Turing machine in time  $O(T(n)^2)$ , so polynomial-time Turing reductions don't actually need to be described using Turing machines. We'll simply take that step on faith.

Complexity-theory researchers prefer to define NP-hardness in terms of polynomial-time **many-one reductions**, which are also called *Karp reductions*. Karp reductions are defined over *languages*: sets of strings over a fixed alphabet  $\Sigma$ . (Without loss of generality, we can assume that  $\Sigma = \{0, 1\}$ .) A many-one reduction from one language  $\Pi' \subseteq \Sigma^*$  to another language  $\Pi \subseteq \Sigma^*$  is a function  $f : \Sigma^* \rightarrow \Sigma^*$  such that  $x \in \Pi'$  if and only if  $f(x) \in \Pi$ . Then we could define a *language*  $\Pi$  to be NP-hard if and only if, for any language  $\Pi' \in \text{NP}$ , there is a many-one reduction from  $\Pi'$  to  $\Pi$  that can be computed in polynomial time.

Every Karp reduction is a Cook reduction, but not vice versa; specifically, any Karp reduction from  $\Pi$  to  $\Pi'$  is equivalent to transforming the input to  $\Pi$  into the input for  $\Pi'$ , invoking an oracle (that is, a subroutine) for  $\Pi'$ , and then returning the answer verbatim. However, as far as we know, not every Cook reduction can be simulated by a Karp reduction.

Complexity theorists prefer Karp reductions primarily because NP is closed under Karp reductions, but is *not* closed under Cook reductions (unless  $\text{NP}=\text{co-NP}$ , which is considered unlikely). There are natural problems that are (1) NP-hard with respect to Cook reductions, but (2) NP-hard with respect to Karp reductions only if  $\text{P}=\text{NP}$ . As a trivial example, consider the problem **UNSAT**: Given a boolean formula, is it *always false*? On the other hand, many-one reductions apply *only* to decision problems (or more formally, to languages); formally, no optimization or construction problem is Karp-NP-hard.

To make things even more confusing, both Cook and Karp originally defined NP-hardness in terms of **logarithmic-space** reductions. Every logarithmic-space reduction is a polynomial-time reduction, but (we think) not vice versa. It is an open question whether relaxing the set of allowed (Cook or Karp) reductions from logarithmic-space to polynomial-time changes the set of NP-hard problems.

Fortunately, none of these subtleties raise their ugly heads in practice, so you can wake up now.

## 29.5 Reductions and SAT

To prove that a problem is NP-hard, we use a *reduction argument*. Reducing problem A to another problem B means describing an algorithm to solve problem A under the assumption that an algorithm for problem B already exists. You're already used to doing reductions, only you probably call it something else, like writing subroutines or utility functions, or modular programming. To prove something is NP-hard, we describe a similar transformation between problems, but not in the direction that most people expect.

You should tattoo the following rule of onto the back of your hand, right next to your Mom's birthday and the *actual* rules of Monopoly.<sup>6</sup>

**To prove that problem A is NP-hard, reduce a known NP-hard problem to A.**

<sup>5</sup>Random-access machines are a model of computation that more faithfully models physical computers. A random-access machine has unbounded random-access memory, modeled as an array  $M[0..\infty]$  where each address  $M[i]$  holds a single  $w$ -bit integer, for some fixed integer  $w$ , and can read to or write from any memory addresses in constant time. RAM algorithms are formally written in assembly-like language, using instructions like **ADD  $i, j, k$**  (meaning " $M[i] \leftarrow M[j] + M[k]$ "), **INDIR  $i, j$**  (meaning " $M[i] \leftarrow M[M[j]]$ "), and **IFZGOTO  $i, l$**  (meaning "if  $M[i] = 0$ , go to line  $l$ "). But in practice, RAM algorithms can be faithfully described in higher-level pseudocode, as long as we're careful about arithmetic precision.

<sup>6</sup>If a player lands on an available property and declines (or is unable) to buy it, that property is immediately auctioned off to the highest bidder; the player who originally declined the property may bid, and bids may be arbitrarily higher or lower than the list price. Players still collect rent while in Jail, but not if the property is mortgaged. A player landing on Free Parking does not win anything. A player landing on Go gets \$200, no more. Finally, Jeff *always* gets the car.

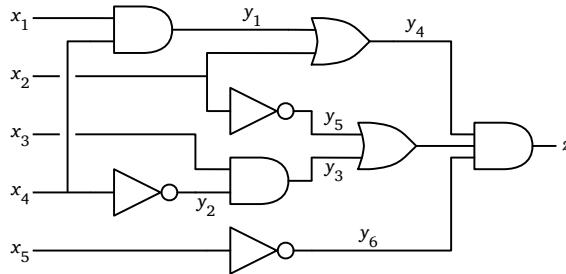
In other words, to prove that your problem is hard, you need to describe an algorithm to solve a *different* problem, which you already know is hard, using a mythical algorithm for *your* problem as a subroutine. The essential logic is a proof by contradiction. Your reduction shows implies that if your problem were easy, then the other problem would be easy, too. Equivalently, since you know the other problem is hard, your problem must also be hard.

For example, consider the *formula satisfiability* problem, usually just called **SAT**. The input to SAT is a boolean *formula* like

$$(a \vee b \vee c \vee \bar{d}) \Leftrightarrow ((b \wedge \bar{c}) \vee (\bar{a} \Rightarrow \bar{d}) \vee (c \neq a \wedge b)),$$

and the question is whether it is possible to assign boolean values to the variables  $a, b, c, \dots$  so that the formula evaluates to TRUE.

To show that SAT is NP-hard, we need to give a reduction from a known NP-hard problem. The only problem we know is NP-hard so far is circuit satisfiability, so let's start there. Given a boolean circuit, we can transform it into a boolean formula by creating new output variables for each gate, and then just writing down the list of gates separated by ANDS. For example, we can transform the example circuit into a formula as follows:

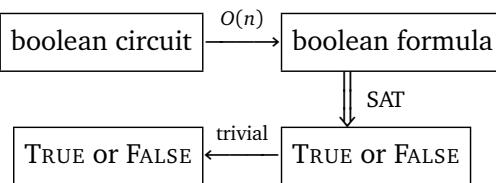


$$\begin{aligned} (y_1 = x_1 \wedge x_4) \wedge (y_2 = \bar{x}_4) \wedge (y_3 = x_3 \wedge y_2) \wedge (y_4 = y_1 \vee x_2) \wedge \\ (y_5 = \bar{x}_2) \wedge (y_6 = \bar{x}_5) \wedge (y_7 = y_3 \vee y_5) \wedge (z = y_4 \wedge y_7 \wedge y_6) \wedge z \end{aligned}$$

A boolean circuit with gate variables added, and an equivalent boolean formula.

Now the original circuit is satisfiable if and only if the resulting formula is satisfiable. Given a satisfying input to the circuit, we can get a satisfying assignment for the formula by computing the output of every gate. Given a satisfying assignment for the formula, we can get a satisfying input to the circuit by just ignoring the internal gate variables  $y_i$  and the output variable  $z$ .

We can transform any boolean circuit into a formula in linear time using depth-first search, and the size of the resulting formula is only a constant factor larger than the size of the circuit. Thus, we have a polynomial-time reduction from circuit satisfiability to SAT:



$$T_{\text{CSAT}}(n) \leq O(n) + T_{\text{SAT}}(O(n)) \implies T_{\text{SAT}}(n) \geq T_{\text{CSAT}}(\Omega(n)) - O(n)$$

The reduction implies that if we had a polynomial-time algorithm for SAT, then we'd have a polynomial-time algorithm for circuit satisfiability, which would imply that P=NP. So SAT is NP-hard.

To prove that a boolean formula is satisfiable, we only have to specify an assignment to the variables that makes the formula TRUE. We can check the proof in linear time just by reading the formula from left to right, evaluating as we go. So SAT is also in NP, and thus is actually NP-complete.

## 29.6 3SAT (from SAT)

A special case of SAT that is particularly useful in proving NP-hardness results is called 3SAT.

A boolean formula is in *conjunctive normal form* (CNF) if it is a conjunction (AND) of several *clauses*, each of which is the disjunction (OR) of several *literal*s, each of which is either a variable or its negation. For example:

$$\overbrace{(a \vee b \vee c \vee d) \wedge (b \vee \bar{c} \vee \bar{d}) \wedge (\bar{a} \vee c \vee d) \wedge (a \vee \bar{b})}^{\text{clause}}$$

A 3CNF formula is a CNF formula with exactly three literals per clause; the previous example is not a 3CNF formula, since its first clause has four literals and its last clause has only two. 3SAT is just SAT restricted to 3CNF formulas: Given a 3CNF formula, is there an assignment to the variables that makes the formula evaluate to TRUE?

We could prove that 3SAT is NP-hard by a reduction from the more general SAT problem, but it's easier just to start over from scratch, with a boolean circuit. We perform the reduction in several stages.

1. *Make sure every AND and OR gate has only two inputs.* If any gate has  $k > 2$  inputs, replace it with a binary tree of  $k - 1$  two-input gates.
2. *Write down the circuit as a formula, with one clause per gate.* This is just the previous reduction.
3. *Change every gate clause into a CNF formula.* There are only three types of clauses, one for each type of gate:

$$\begin{aligned} a = b \wedge c &\longrightarrow (a \vee \bar{b} \vee \bar{c}) \wedge (\bar{a} \vee b) \wedge (\bar{a} \vee c) \\ a = b \vee c &\longrightarrow (\bar{a} \vee b \vee c) \wedge (a \vee \bar{b}) \wedge (a \vee \bar{c}) \\ a = \bar{b} &\longrightarrow (a \vee b) \wedge (\bar{a} \vee \bar{b}) \end{aligned}$$

4. *Make sure every clause has exactly three literals.* Introduce new variables into each one- and two-literal clause, and expand it into two clauses as follows:

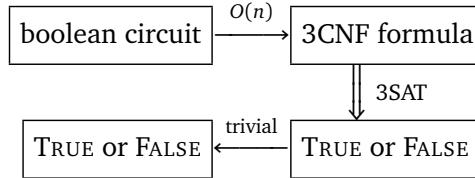
$$\begin{aligned} a &\longrightarrow (a \vee x \vee y) \wedge (a \vee \bar{x} \vee y) \wedge (a \vee x \vee \bar{y}) \wedge (a \vee \bar{x} \vee \bar{y}) \\ a \vee b &\longrightarrow (a \vee b \vee x) \wedge (a \vee b \vee \bar{x}) \end{aligned}$$

For example, if we start with the same example circuit we used earlier, we obtain the following 3CNF formula. Although this may look a lot more ugly and complicated than the original circuit at first glance, it's actually only a constant factor larger—every binary gate in the original circuit has been transformed into at most five clauses. Even if the formula size were a large *polynomial* function (like  $n^{573}$ ) of the

circuit size, we would still have a valid reduction.

$$\begin{aligned}
 & (y_1 \vee \bar{x}_1 \vee \bar{x}_4) \wedge (\bar{y}_1 \vee x_1 \vee z_1) \wedge (\bar{y}_1 \vee x_1 \vee \bar{z}_1) \wedge (\bar{y}_1 \vee x_4 \vee z_2) \wedge (\bar{y}_1 \vee x_4 \vee \bar{z}_2) \\
 & \wedge (y_2 \vee x_4 \vee z_3) \wedge (y_2 \vee x_4 \vee \bar{z}_3) \wedge (\bar{y}_2 \vee \bar{x}_4 \vee z_4) \wedge (\bar{y}_2 \vee \bar{x}_4 \vee \bar{z}_4) \\
 & \wedge (y_3 \vee \bar{x}_3 \vee \bar{y}_2) \wedge (\bar{y}_3 \vee x_3 \vee z_5) \wedge (\bar{y}_3 \vee x_3 \vee \bar{z}_5) \wedge (\bar{y}_3 \vee y_2 \vee z_6) \wedge (\bar{y}_3 \vee y_2 \vee \bar{z}_6) \\
 & \wedge (\bar{y}_4 \vee y_1 \vee x_2) \wedge (y_4 \vee \bar{x}_2 \vee z_7) \wedge (y_4 \vee \bar{x}_2 \vee \bar{z}_7) \wedge (y_4 \vee \bar{y}_1 \vee z_8) \wedge (y_4 \vee \bar{y}_1 \vee \bar{z}_8) \\
 & \wedge (y_5 \vee x_2 \vee z_9) \wedge (y_5 \vee x_2 \vee \bar{z}_9) \wedge (\bar{y}_5 \vee \bar{x}_2 \vee z_{10}) \wedge (\bar{y}_5 \vee \bar{x}_2 \vee \bar{z}_{10}) \\
 & \wedge (y_6 \vee x_5 \vee z_{11}) \wedge (y_6 \vee x_5 \vee \bar{z}_{11}) \wedge (\bar{y}_6 \vee \bar{x}_5 \vee z_{12}) \wedge (\bar{y}_6 \vee \bar{x}_5 \vee \bar{z}_{12}) \\
 & \wedge (\bar{y}_7 \vee y_3 \vee y_5) \wedge (y_7 \vee \bar{y}_3 \vee z_{13}) \wedge (y_7 \vee \bar{y}_3 \vee \bar{z}_{13}) \wedge (y_7 \vee \bar{y}_5 \vee z_{14}) \wedge (y_7 \vee \bar{y}_5 \vee \bar{z}_{14}) \\
 & \wedge (y_8 \vee \bar{y}_4 \vee \bar{y}_7) \wedge (\bar{y}_8 \vee y_4 \vee z_{15}) \wedge (\bar{y}_8 \vee y_4 \vee \bar{z}_{15}) \wedge (\bar{y}_8 \vee y_7 \vee z_{16}) \wedge (\bar{y}_8 \vee y_7 \vee \bar{z}_{16}) \\
 & \wedge (y_9 \vee \bar{y}_8 \vee \bar{y}_6) \wedge (\bar{y}_9 \vee y_8 \vee z_{17}) \wedge (\bar{y}_9 \vee y_8 \vee \bar{z}_{17}) \wedge (\bar{y}_9 \vee y_6 \vee z_{18}) \wedge (\bar{y}_9 \vee y_6 \vee \bar{z}_{18}) \\
 & \wedge (y_9 \vee z_{19} \vee z_{20}) \wedge (y_9 \vee \bar{z}_{19} \vee z_{20}) \wedge (y_9 \vee z_{19} \vee \bar{z}_{20}) \wedge (y_9 \vee \bar{z}_{19} \vee \bar{z}_{20})
 \end{aligned}$$

This process transforms the circuit into an equivalent 3CNF formula; the output formula is satisfiable if and only if the input circuit is satisfiable. As with the more general SAT problem, the formula is only a constant factor larger than any reasonable description of the original circuit, and the reduction can be carried out in polynomial time. Thus, we have a polynomial-time reduction from circuit satisfiability to 3SAT:



$$T_{\text{CSAT}}(n) \leq O(n) + T_{\text{3SAT}}(O(n)) \implies T_{\text{3SAT}}(n) \geq T_{\text{CSAT}}(\Omega(n)) - O(n)$$

We conclude 3SAT is NP-hard. And because 3SAT is a special case of SAT, it is also in NP. Therefore, 3SAT is NP-complete.

## 29.7 Maximum Independent Set (from 3SAT)

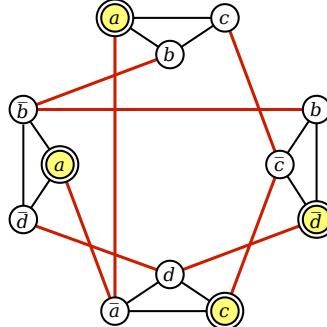
For the next few problems we consider, the input is a simple, unweighted graph, and the problem asks for the size of the largest or smallest subgraph satisfying some structural property.

Let  $G$  be an arbitrary graph. An **independent set** in  $G$  is a subset of the vertices of  $G$  with no edges between them. The **maximum independent set** problem, or simply **MAXINDSET**, asks for the size of the largest independent set in a given graph.

I'll prove that MAXINDSET is NP-hard (but not NP-complete, since it isn't a decision problem) using a reduction from 3SAT. I'll describe a reduction from a 3CNF formula into a graph that has an independent set of a certain size if and only if the formula is satisfiable. The graph has one node for each instance of each literal in the formula. Two nodes are connected by an edge if (1) they correspond to literals in the same clause, or (2) they correspond to a variable and its inverse. For example, the formula  $(a \vee b \vee c) \wedge (b \vee \bar{c} \vee \bar{d}) \wedge (\bar{a} \vee c \vee d) \wedge (a \vee \bar{b} \vee \bar{d})$  is transformed into the following graph.

Now suppose the original formula had  $k$  clauses. Then I claim that the formula is satisfiable if and only if the graph has an independent set of size  $k$ .

1. **independent set  $\implies$  satisfying assignment:** If the graph has an independent set of  $k$  vertices, then each vertex must come from a different clause. To obtain a satisfying assignment, we assign

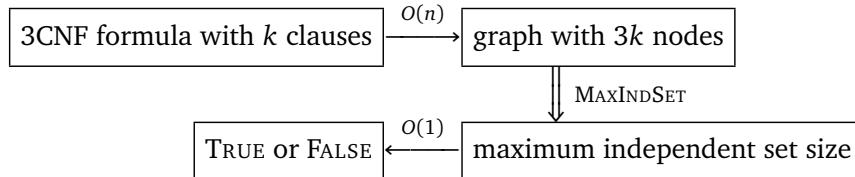


A graph derived from a 3CNF formula, and an independent set of size 4.  
Black edges join literals from the same clause; red (heavier) edges join contradictory literals.

the value TRUE to each literal in the independent set. Since contradictory literals are connected by edges, this assignment is consistent. There may be variables that have no literal in the independent set; we can set these to any value we like. The resulting assignment satisfies the original 3CNF formula.

2. **satisfying assignment  $\implies$  independent set:** If we have a satisfying assignment, then we can choose one literal in each clause that is TRUE. Those literals form an independent set in the graph.

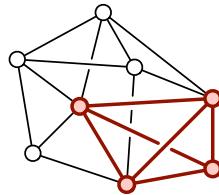
Thus, the reduction is correct. Since the reduction from 3CNF formula to graph takes polynomial time, we conclude that MAXINDSET is NP-hard. Here's a diagram of the reduction:



$$T_{\text{3SAT}}(n) \leq O(n) + T_{\text{MAXINDSET}}(O(n)) \implies T_{\text{MAXINDSET}}(n) \geq T_{\text{3SAT}}(\Omega(n)) - O(n)$$

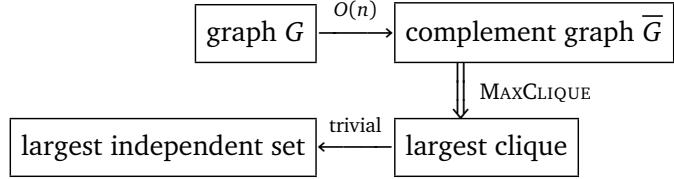
## 29.8 Clique (from Independent Set)

A *clique* is another name for a complete graph, that is, a graph where every pair of vertices is connected by an edge. The *maximum clique size* problem, or simply MAXCLIQUE, is to compute, given a graph, the number of nodes in its largest complete subgraph.



A graph with maximum clique size 4.

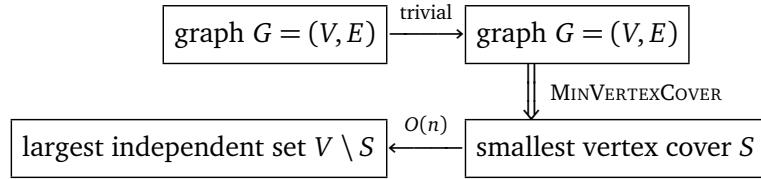
There is an easy proof that MAXCLIQUE is NP-hard, using a reduction from MAXINDSET. Any graph  $G$  has an *edge-complement*  $\bar{G}$  with the same vertices, but with exactly the opposite set of edges— $(u, v)$  is an edge in  $\bar{G}$  if and only if it is *not* an edge in  $G$ . A set of vertices is independent in  $G$  if and only if the same vertices define a clique in  $\bar{G}$ . Thus, we can compute the largest independent in a graph simply by computing the largest clique in the complement of the graph.



## 29.9 Vertex Cover (from Independent Set)

A *vertex cover* of a graph is a set of vertices that touches every edge in the graph. The *MINVERTEXCOVER* problem is to find the smallest vertex cover in a given graph.

Again, the proof of NP-hardness is simple, and relies on just one fact: If  $I$  is an independent set in a graph  $G = (V, E)$ , then  $V \setminus I$  is a vertex cover. Thus, to find the *largest* independent set, we just need to find the vertices that aren't in the *smallest* vertex cover of the same graph.

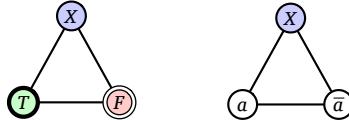


## 29.10 Graph Coloring (from 3SAT)

A  $k$ -*coloring* of a graph is a map  $C: V \rightarrow \{1, 2, \dots, k\}$  that assigns one of  $k$  ‘colors’ to each vertex, so that every edge has two different colors at its endpoints. The graph coloring problem is to find the smallest possible number of colors in a legal coloring. To show that this problem is NP-hard, it’s enough to consider the special case *3COLORABLE*: Given a graph, does it have a 3-coloring?

To prove that *3COLORABLE* is NP-hard, we use a reduction from *3SAT*. Given a 3CNF formula  $\Phi$ , we produce a graph  $G_\Phi$  as follows. The graph consists of a *truth* gadget, one *variable* gadget for each variable in the formula, and one *clause* gadget for each clause in the formula.

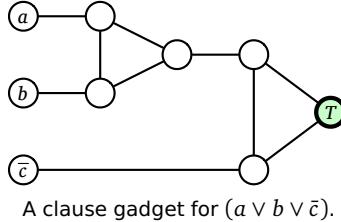
- The truth gadget is just a triangle with three vertices  $T$ ,  $F$ , and  $X$ , which intuitively stand for TRUE, FALSE, and OTHER. Since these vertices are all connected, they must have different colors in any 3-coloring. For the sake of convenience, we will *name* those colors TRUE, FALSE, and OTHER. Thus, when we say that a node is colored TRUE, all we mean is that it must be colored the same as the node  $T$ .
- The variable gadget for a variable  $a$  is also a triangle joining two new nodes labeled  $a$  and  $\bar{a}$  to node  $X$  in the truth gadget. Node  $a$  must be colored either TRUE or FALSE, and so node  $\bar{a}$  must be colored either FALSE or TRUE, respectively.



The truth gadget and a variable gadget for  $a$ .

- Finally, each clause gadget joins three literal nodes to node  $T$  in the truth gadget using five new unlabeled nodes and ten edges; see the figure below. A straightforward case analysis implies that if all three literal nodes in the clause gadget are colored FALSE, then some edge in the gadget must be monochromatic. Since the variable gadgets force each literal node to be colored either TRUE or

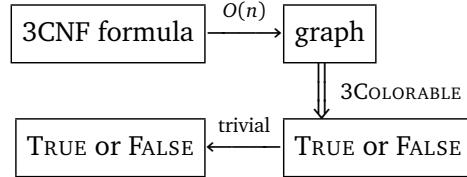
FALSE, in any valid 3-coloring, at least one of the three literal nodes is colored TRUE. On the other hand, for any coloring of the literal nodes where at least one literal node is colored TRUE, there is a valid 3-coloring of the clause gadget.



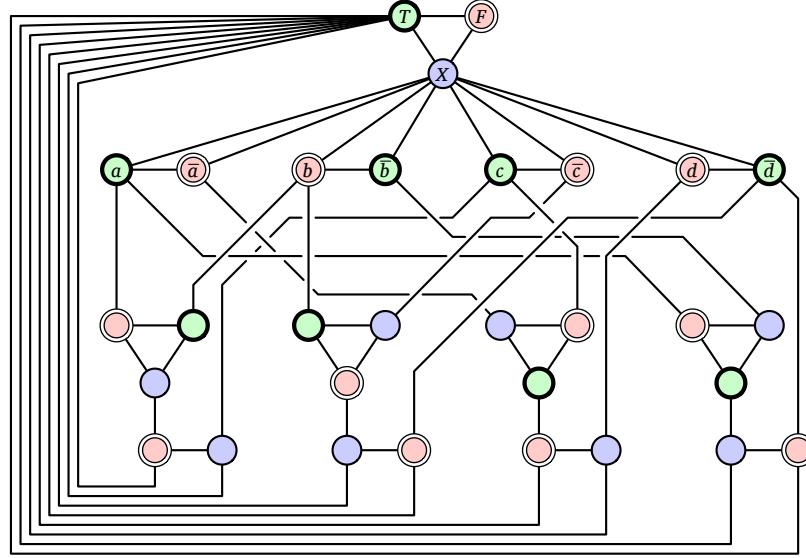
A clause gadget for  $(a \vee b \vee \bar{c})$ .

I need to emphasize here that the final graph  $G_\Phi$  contains exactly *one* node  $T$ , exactly *one* node  $F$ , and exactly *two* nodes  $a$  and  $\bar{a}$  for each variable.

Now the proof of correctness is just brute force case analysis. If the graph is 3-colorable, then we can extract a satisfying assignment from any 3-coloring—at least one of the three literal nodes in every clause gadget is colored TRUE. Conversely, if the formula is satisfiable, then we can color the graph according to any satisfying assignment.



For example, the formula  $(a \vee b \vee c) \wedge (b \vee \bar{c} \vee \bar{d}) \wedge (\bar{a} \vee c \vee d) \wedge (a \vee \bar{b} \vee \bar{d})$  that I used to illustrate the MAXCLIQUE reduction would be transformed into the following graph. The 3-coloring is one of several that correspond to the satisfying assignment  $a = c = \text{TRUE}$ ,  $b = d = \text{FALSE}$ .



A 3-colorable graph derived from a satisfiable 3CNF formula.

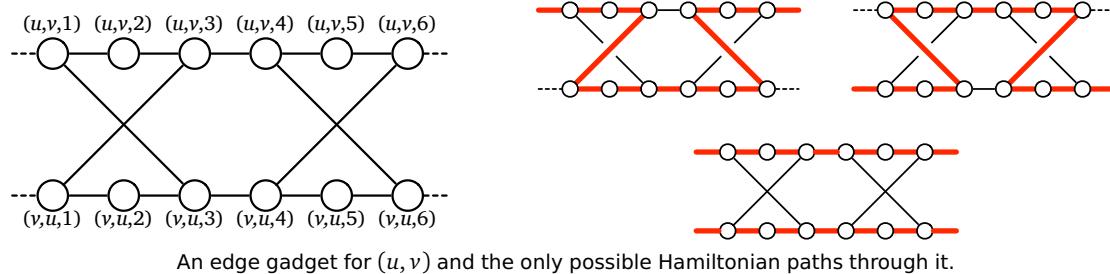
We can easily verify that a graph has been correctly 3-colored in linear time: just compare the endpoints of every edge. Thus, 3COLORING is in NP, and therefore NP-complete. Moreover, since 3COLORING is a special case of the more general graph coloring problem—What is the minimum number of colors?—the more problem is also NP-hard, but *not* NP-complete, because it's not a decision problem.

### 29.11 Hamiltonian Cycle (from Vertex Cover)

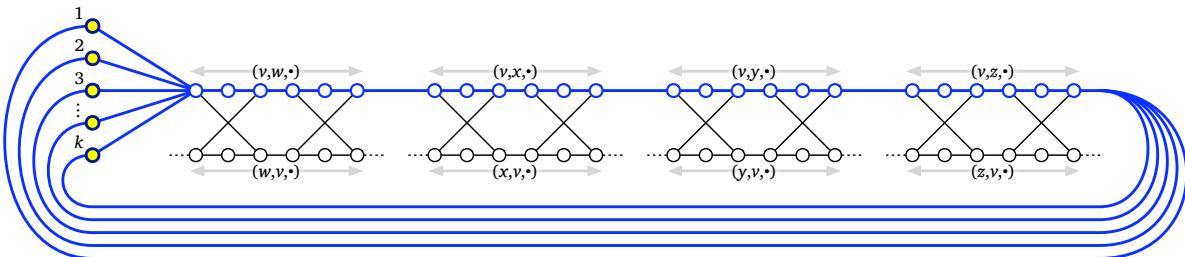
A *Hamiltonian cycle* in a graph is a cycle that visits every vertex exactly once. This is very different from an *Eulerian cycle*, which is actually a closed *walk* that traverses every *edge* exactly once. Eulerian cycles are easy to find and construct in linear time using a variant of depth-first search. Determining whether a graph contains a Hamiltonian cycle, on the other hand, is NP-hard.

To prove this, we describe a reduction from the vertex cover problem. Given a graph  $G$  and an integer  $k$ , we need to transform it into another graph  $G'$ , such that  $G'$  has a Hamiltonian cycle if and only if  $G$  has a vertex cover of size  $k$ . As usual, our transformation uses several gadgets.

- For each edge  $uv$  in  $G$ , we have an *edge gadget* in  $G'$  consisting of twelve vertices and fourteen edges, as shown below. The four corner vertices  $(u, v, 1), (u, v, 6), (v, u, 1)$ , and  $(v, u, 6)$  each have an edge leaving the gadget. A Hamiltonian cycle can only pass through an edge gadget in only three ways. Eventually, these options will correspond to one or both vertices  $u$  and  $v$  being in the vertex cover.



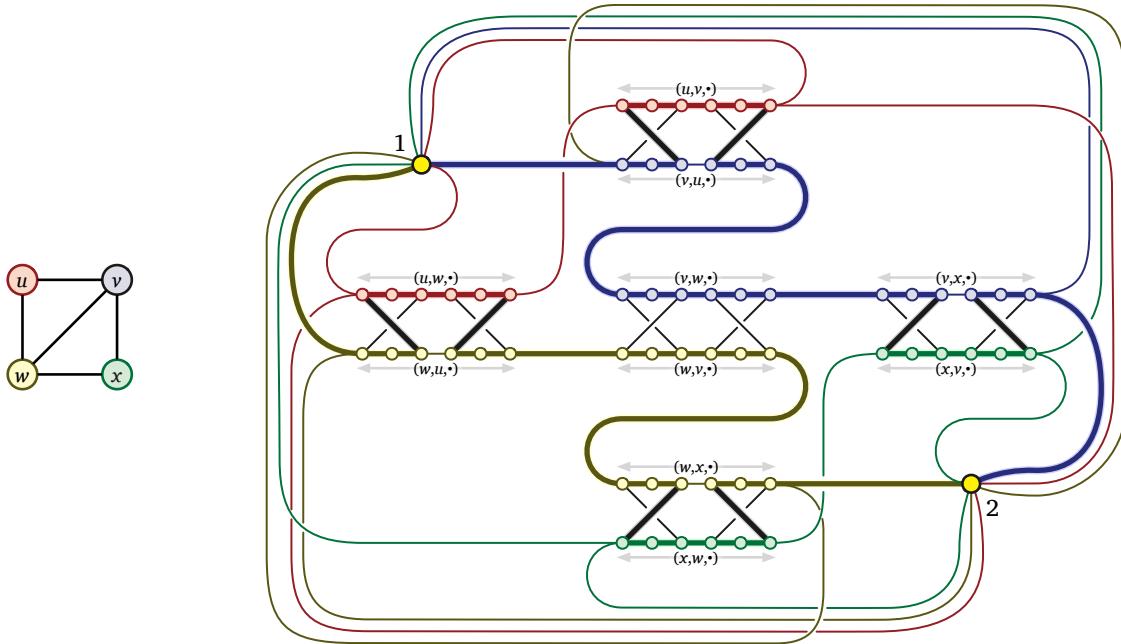
- $G'$  also contains  $k$  *cover vertices*, simply numbered 1 through  $k$ .
- Finally, for each vertex  $u$  in  $G$ , we string together all the edge gadgets for edges  $(u, v)$  into a single *vertex chain*, and then connect the ends of the chain to all the cover vertices. Specifically, suppose vertex  $u$  has  $d$  neighbors  $v_1, v_2, \dots, v_d$ . Then  $G'$  has  $d - 1$  edges between  $(u, v_i, 6)$  and  $(u, v_{i+1}, 1)$ , plus  $k$  edges between the cover vertices and  $(u, v_1, 1)$ , and finally  $k$  edges between the cover vertices and  $(u, v_d, 6)$ .



The vertex chain for  $v$ : all edge gadgets involving  $v$  are strung together and joined to the  $k$  cover vertices.

An example of the output from our reduction is shown on the next page.

It is now straightforward but tedious to prove that if  $\{v_1, v_2, \dots, v_k\}$  is a vertex cover of  $G$ , then  $G'$  has a Hamiltonian cycle—start at cover vertex 1, through traverse the vertex chain for  $v_1$ , then visit cover vertex 2, then traverse the vertex chain for  $v_2$ , and so forth, eventually returning to cover vertex 1. Conversely, any Hamiltonian cycle in  $G'$  alternates between cover vertices and vertex chains, and the vertex chains correspond to the  $k$  vertices in a vertex cover of  $G$ . (This is a little harder to prove.) Thus,  $G$  has a vertex cover of size  $k$  if and only if  $G'$  has a Hamiltonian cycle.



The original graph  $G$  with vertex cover  $\{v, w\}$ , and the transformed graph  $G'$  with a corresponding Hamiltonian cycle.  
Vertex chains are colored to match their corresponding vertices.

The transformation from  $G$  to  $G'$  takes at most  $O(n^2)$  time; we conclude that the Hamiltonian cycle problem is NP-hard. Moreover, since we can easily verify a Hamiltonian cycle in linear time, the Hamiltonian cycle problem is in NP, and therefore is NP-complete.

A closely related problem to Hamiltonian cycles is the famous *traveling salesman problem*—Given a weighted graph  $G$ , find the shortest cycle that visits every vertex. Finding the shortest cycle is obviously harder than determining if a cycle exists at all, so the traveling salesman problem is also NP-hard.

## 29.12 Subset Sum (from Vertex Cover)

The last problem that we will prove NP-hard is the **SUBSETSUM** problem considered in the very first lecture on recursion: Given a set  $X$  of positive integers and an integer  $t$ , determine whether  $X$  has a subset whose elements sum to  $t$ .

To prove this problem is NP-hard, we apply another reduction from the vertex cover problem. Given a graph  $G$  and an integer  $k$ , we need to transform it into set of integers  $X$  and an integer  $t$ , such that  $X$  has a subset that sums to  $t$  if and only if  $G$  has a vertex cover of size  $k$ . Our transformation uses just two ‘gadgets’; these are *integers* representing vertices and edges in  $G$ .

Number the *edges* of  $G$  arbitrarily from 0 to  $m - 1$ . Our set  $X$  contains the integer  $b_i := 4^i$  for each edge  $i$ , and the integer

$$a_v := 4^m + \sum_{i \in \Delta(v)} 4^i$$

for each vertex  $v$ , where  $\Delta(v)$  is the set of edges that have  $v$  as an endpoint. Alternately, we can think of each integer in  $X$  as an  $(m + 1)$ -digit number written in base 4. The  $m$ th digit is 1 if the integer represents a vertex, and 0 otherwise; and for each  $i < m$ , the  $i$ th digit is 1 if the integer represents edge  $i$  or one of its endpoints, and 0 otherwise. Finally, we set the target sum

$$t := k \cdot 4^m + \sum_{i=0}^{m-1} 2 \cdot 4^i.$$

Now let's prove that the reduction is correct. First, suppose there is a vertex cover of size  $k$  in the original graph  $G$ . Consider the subset  $X_C \subseteq X$  that includes  $a_v$  for every vertex  $v$  in the vertex cover, and  $b_i$  for every edge  $i$  that has *exactly one* vertex in the cover. The sum of these integers, written in base 4, has a 2 in each of the first  $m$  digits; in the most significant digit, we are summing exactly  $k$  1's. Thus, the sum of the elements of  $X_C$  is exactly  $t$ .

On the other hand, suppose there is a subset  $X' \subseteq X$  that sums to  $t$ . Specifically, we must have

$$\sum_{v \in V'} a_v + \sum_{i \in E'} b_i = t$$

for some subsets  $V' \subseteq V$  and  $E' \subseteq E$ . Again, if we sum these base-4 numbers, there are no carries in the first  $m$  digits, because for each  $i$  there are only three numbers in  $X$  whose  $i$ th digit is 1. Each edge number  $b_i$  contributes only one 1 to the  $i$ th digit of the sum, but the  $i$ th digit of  $t$  is 2. Thus, for each edge in  $G$ , at least one of its endpoints must be in  $V'$ . In other words,  $V'$  is a vertex cover. On the other hand, only vertex numbers are larger than  $4^m$ , and  $\lfloor t/4^m \rfloor = k$ , so  $V'$  has at most  $k$  elements. (In fact, it's not hard to see that  $V'$  has *exactly*  $k$  elements.)

For example, given the four-vertex graph used on the previous page to illustrate the reduction to Hamiltonian cycle, our set  $X$  might contain the following base-4 integers:

$$\begin{array}{ll} a_u := 111000_4 = 1344 & b_{uv} := 010000_4 = 256 \\ a_v := 110110_4 = 1300 & b_{uw} := 001000_4 = 64 \\ a_w := 101101_4 = 1105 & b_{vw} := 000100_4 = 16 \\ a_x := 100011_4 = 1029 & b_{vx} := 000010_4 = 4 \\ & b_{wx} := 000001_4 = 1 \end{array}$$

If we are looking for a vertex cover of size 2, our target sum would be  $t := 222222_4 = 2730$ . Indeed, the vertex cover  $\{v, w\}$  corresponds to the subset  $\{a_v, a_w, b_{uv}, b_{uw}, b_{vx}, b_{wx}\}$ , whose sum is  $1300 + 1105 + 256 + 64 + 4 + 1 = 2730$ .

The reduction can clearly be performed in polynomial time. Since VERTEXCOVER is NP-hard, it follows that SUBSETSUM is NP-hard.

There is one subtle point that needs to be emphasized here. Way back at the beginning of the semester, we developed a dynamic programming algorithm to solve SUBSETSUM in time  $O(nt)$ . Isn't this a polynomial-time algorithm? idn't we just prove that P=NP? Hey, where's our million dollars? Alas, life is not so simple. True, the running time is polynomial in  $n$  and  $t$ , but in order to qualify as a true polynomial-time algorithm, the running time must be a polynomial function of the size of the input. The values of the elements of  $X$  and the target sum  $t$  could be exponentially larger than the number of input bits. Indeed, the reduction we just described produces a value of  $t$  that is exponentially larger than the size of our original input graph, which would force our dynamic programming algorithm to run in exponential time.

Algorithms like this are said to run in *pseudo-polynomial time*, and any NP-hard problem with such an algorithm is called **weakly NP-hard**. Equivalently, a weakly NP-hard problem is one that can be solved in polynomial time when all input numbers are represented in *unary* (as a sum of 1s), but becomes NP-hard when all input numbers are represented in *binary*. If a problem is NP-hard even when all the input numbers are represented in unary, we say that the problem is **strongly NP-hard**.

### 29.13 Other Useful NP-hard Problems

Literally thousands of different problems have been proved to be NP-hard. I want to close this note by listing a few NP-hard problems that are useful in deriving reductions. I won't describe the NP-hardness

proofs for these problems in detail, but you can find most of them in Garey and Johnson's classic *Scary Black Book of NP-Completeness*.<sup>7</sup>

- PLANARCIRCUITSAT: Given a boolean circuit that can be embedded in the plane so that no two wires cross, is there an input that makes the circuit output TRUE? This problem can be proved NP-hard by reduction from the general circuit satisfiability problem, by replacing each crossing with a small series of gates. (The details of the reduction are an easy<sup>8</sup> exercise.)
- NOTALLEQUAL3SAT: Given a 3CNF formula, is there an assignment of values to the variables so that every clause contains at least one TRUE literal *and* at least one FALSE literal? This problem can be proved NP-hard by reduction from the usual 3SAT.
- PLANAR3SAT: Given a 3CNF boolean formula, consider a bipartite graph whose vertices are the clauses and variables, where an edge indicates that a variable (or its negation) appears in a clause. If this graph is planar, the 3CNF formula is also called planar. The PLANAR3SAT problem asks, given a planar 3CNF formula, whether it has a satisfying assignment. This problem can be proved NP-hard by reduction from PLANARCIRCUITSAT.<sup>9</sup>
- EXACT3DIMENSIONALMATCHING or X3M: Given a set  $S$  and a collection of three-element subsets of  $S$ , called *triples*, is there a sub-collection of disjoint triples that exactly cover  $S$ ? This problem can be proved NP-hard by a reduction from 3SAT.
- PARTITION: Given a set  $S$  of  $n$  integers, are there subsets  $A$  and  $B$  such that  $A \cup B = S$ ,  $A \cap B = \emptyset$ , and

$$\sum_{a \in A} a = \sum_{b \in B} b?$$

This problem can be proved NP-hard by a simple reduction from SUBSETSUM. Like SUBSETSUM, the PARTITION problem is only weakly NP-hard.

- 3PARTITION: Given a set  $S$  of  $3n$  integers, can it be partitioned into  $n$  disjoint three-element subsets, such that every subset has exactly the same sum? Despite the similar names, this problem is *very* different from PARTITION; sorry, I didn't make up the names. This problem can be proved NP-hard by reduction from X3M. Unlike PARTITION, the 3PARTITION problem is *strongly* NP-hard, that is, it remains NP-hard even if the input numbers are less than some polynomial in  $n$ .
- SETCOVER: Given a collection of sets  $\mathcal{S} = \{S_1, S_2, \dots, S_m\}$ , find the smallest sub-collection of  $S_i$ 's that contains all the elements of  $\bigcup_i S_i$ . This problem is a generalization of both VERTEXCOVER and X3M.
- HITTINGSET: Given a collection of sets  $\mathcal{S} = \{S_1, S_2, \dots, S_m\}$ , find the minimum number of elements of  $\bigcup_i S_i$  that hit every set in  $\mathcal{S}$ . This problem is also a generalization of VERTEXCOVER.
- LONGESTPATH: Given a non-negatively weighted graph  $G$  and two vertices  $u$  and  $v$ , what is the longest simple path from  $u$  to  $v$  in the graph? A path is *simple* if it visits each vertex at most once. This problem is a generalization of the HAMILTONIANPATH problem. Of course, the corresponding *shortest* path problem is in P.

<sup>7</sup>Michael Garey and David Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman and Co., 1979.

<sup>8</sup>or at least *nondeterministically* easy

<sup>9</sup>Surprisingly, PLANARNOTALLEQUAL3SAT is solvable in polynomial time!

- STEINERTREE: Given a weighted, undirected graph  $G$  with some of the vertices marked, what is the minimum-weight subtree of  $G$  that contains every marked vertex? If every vertex is marked, the minimum Steiner tree is just the minimum spanning tree; if exactly two vertices are marked, the minimum Steiner tree is just the shortest path between them. This problem can be proved NP-hard by reduction to HAMILTONIANPATH.

In addition to these dry but useful problems, most interesting puzzles and solitaire games have been shown to be NP-hard, or to have NP-hard generalizations. (Arguably, if a game or puzzle isn't at least NP-hard, it isn't interesting!) Some familiar examples include Minesweeper (by reduction from CIRCUITSAT)<sup>10</sup>, Tetris (by reduction from 3PARTITION)<sup>11</sup>, and Sudoku (by reduction from 3SAT)<sup>12</sup>.

#### \*29.14 On Beyond Zebra

P and NP are only the first two steps in an enormous hierarchy of complexity classes. To close these notes, let me describe a few more classes of interest.

**Polynomial Space.** **PSPACE** is the set of decision problems that can be solved using polynomial space. Every problem in NP (and therefore in P) is also in PSPACE. It is generally believed that  $NP \neq PSPACE$ , but nobody can even prove that  $P \neq PSPACE$ . A problem  $\Pi$  is **PSPACE-hard** if, for any problem  $\Pi'$  that can be solved using polynomial space, there is a polynomial-time many-one reduction from  $\Pi'$  to  $\Pi$ . A problem is **PSPACE-complete** if it is both PSPACE-hard and in PSPACE. If any PSPACE-hard problem is in NP, then  $PSPACE=NP$ ; similarly, if any PSPACE-hard problem is in P, then  $PSPACE=P$ .

The canonical PSPACE-complete problem is the *quantified boolean formula* problem, or **QBF**: Given a boolean formula  $\Phi$  that may include any number of universal or existential quantifiers, but no free variables, is  $\Phi$  equivalent to TRUE? For example, the following expression is a valid input to QBF:

$$\exists a : \forall b : \exists c : (\forall d : a \vee b \vee c \vee \bar{d}) \Leftrightarrow ((b \wedge \bar{c}) \vee (\exists e : (\bar{a} \Rightarrow e) \vee (c \neq a \wedge e))).$$

SAT is provably equivalent the special case of QBF where the input formula contains only existential quantifiers. QBF remains PSPACE-hard even when the input formula must have all its quantifiers at the beginning, the quantifiers strictly alternate between  $\exists$  and  $\forall$ , and the quantified proposition is in conjunctive normal form, with exactly three literals in each clause, for example:

$$\exists a : \forall b : \exists c : \forall d : ((a \vee b \vee c) \wedge (b \vee \bar{c} \vee \bar{d}) \wedge (\bar{a} \vee c \vee d) \wedge (a \vee \bar{b} \vee \bar{d}))$$

This restricted version of QBF can also be phrased as a two-player strategy question. Suppose two players, Alice and Bob, are given a 3CNF predicate with free variables  $x_1, x_2, \dots, x_n$ . The players alternately assign values to the variables in order by index—Alice assigns a value to  $x_1$ , Bob assigns a value to  $x_2$ , and so on. Alice eventually assigns values to every variable with an odd index, and Bob eventually assigns values to every variable with an even index. Alice wants to make the expression TRUE, and Bob wants to make it FALSE. Assuming Alice and Bob play perfectly, who wins this game?

<sup>10</sup>Richard Kaye. Minesweeper is NP-complete. *Mathematical Intelligencer* 22(2):9–15, 2000. <http://www.mat.bham.ac.uk/R.W.Kaye/minesw/minesw.pdf>

<sup>11</sup>Ron Breukelaar\*, Erik D. Demaine, Susan Hohenberger\*, Hendrik J. Hoogeboom, Walter A. Kosters, and David Liben-Nowell\*. Tetris is Hard, Even to Approximate. *International Journal of Computational Geometry and Applications* 14:41–68, 2004. The reduction was considerably simplified between its discovery in 2002 and its publication in 2004.

<sup>12</sup>Takayuki Yato and Takahiro Seta. Complexity and Completeness of Finding Another Solution and Its Application to Puzzles. *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences* E86-A(5):1052–1060, 2003. <http://www-imai.is.s.u-tokyo.ac.jp/~yato/data2/MasterThesis.pdf>.

Not surprisingly, most two-player games<sup>13</sup> like tic-tac-toe, reversi, checkers, go, chess, and mancala—or more accurately, appropriate generalizations of these constant-size games to arbitrary board sizes—are PSPACE-hard.

Another canonical PSPACE-hard problem is *NFA totality*: Given a non-deterministic finite-state automaton  $M$  over some alphabet  $\Sigma$ , does  $M$  accept every string in  $\Sigma^*$ ? The closely related problems *NFA equivalence* (Do two given NFAs accept the same language?) and *NFA minimization* (Find the smallest NFA that accepts the same language as a given NFA) are also PSPACE-hard, as are the corresponding questions about regular expressions. (The corresponding questions about *deterministic* finite-state automata are all solvable in polynomial time.)

**Exponential time.** The next significantly larger complexity class, **EXP** (also called EXPTIME), is the set of decision problems that can be solved in exponential time, that is, using at most  $2^{n^c}$  steps for some constant  $c > 0$ . Every problem in PSPACE (and therefore in NP (and therefore in P)) is also in EXP. It is generally believed that PSPACE  $\subsetneq$  EXP, but nobody can even prove that  $\text{NP} \neq \text{EXP}$ . A problem  $\Pi$  is **EXP-hard** if, for any problem  $\Pi'$  that can be solved in *exponential* time, there is a *polynomial*-time many-one reduction from  $\Pi'$  to  $\Pi$ . A problem is **EXP-complete** if it is both EXP-hard and in EXP. If any EXP-hard problem is in PSPACE, then EXP=PSPACE; similarly, if any EXP-hard problem is in NP, then EXP=NP. We do know that  $\text{P} \neq \text{EXP}$ ; in particular, no EXP-hard problem is in P.

Natural generalizations of many interesting 2-player games—like checkers, chess, mancala, and go—are actually EXP-hard. The boundary between PSPACE-complete games and EXP-hard games is rather subtle. For example, there are three ways to draw in chess (the standard  $8 \times 8$  game): stalemate (the player to move is not in check but has no legal moves), repeating the same board position three times, or moving fifty times without capturing a piece. The  $n \times n$  generalization of chess is either in PSPACE or EXP-hard depending on how we generalize these rules. If we declare a draw after (say)  $n^3$  capture-free moves, then every game must end after a polynomial number of moves, so we can simulate all possible games from any given position using only polynomial space. On the other hand, if we ignore the capture-free move rule entirely, the resulting game can last an exponential number of moves, so there is no obvious way to detect a repeating position using only polynomial space; indeed, this version of  $n \times n$  chess is EXP-hard.

**Excelsior!** Naturally, even exponential time is not the end of the story. **NEXP** is the class of decision problems that can be solved in *nondeterministic* exponential time; equivalently, a decision problem is in NEXP if and only if, for every YES instance, there is a *proof* of this fact that can be checked in exponential time. **EXPSPACE** is the set of decision problems that can be solved using exponential space. Even these larger complexity classes have hard and complete problems; for example, if we add the intersection operator  $\cap$  to the syntax of regular expressions, deciding whether two such expressions describe the same language is EXPSPACE-hard. Beyond EXPSPACE are complexity classes with *doubly*-exponential resource bounds (EEXP, NEXP, and EXPSPACE), then *triply* exponential resource bounds (EEEXP, NEEXP, and EXPSPACE), and so on ad infinitum.

All these complexity classes can be ordered by inclusion as follows:

$$\text{P} \subseteq \text{NP} \subseteq \text{PSPACE} \subseteq \text{EXP} \subseteq \text{NEXP} \subseteq \text{EXPSPACE} \subseteq \text{EEXP} \subseteq \text{NEEXP} \subseteq \text{EXPSPACE} \subseteq \text{EEEXP} \subseteq \dots,$$

Most complexity theorists strongly believe that every inclusion in this sequence is strict; that is, no two of these complexity classes are equal. However, the strongest result that has been proved is that every

---

<sup>13</sup>For a good (but now slightly dated) overview of known results on the computational complexity of games and puzzles, see Erik D. Demaine and Robert Hearn's survey "Playing Games with Algorithms: Algorithmic Combinatorial Game Theory" at <http://arxiv.org/abs/cs.CC/0106019>.

class in this sequence is strictly contained in the class *three* steps later in the sequence. For example, we have proofs that  $P \neq EXP$  and  $PSPACE \neq EXPSPACE$ , but not whether  $P \neq PSPACE$  or  $NP \neq EXP$ .

The limit of this series of increasingly exponential complexity classes is the class **ELEMENTARY** of decision problems that can be solved using time or space bounded by any so-called elementary function of  $n$ . A function of  $n$  is *elementary* if it is equal to  $2 \uparrow^k p(n)$  for some integer  $k$  and some polynomial  $p(n)$ , where

$$2 \uparrow^k x := \begin{cases} x & \text{if } k = 0, \\ 2^{2 \uparrow^{k-1} x} & \text{otherwise.} \end{cases}$$

You might be tempted to conjecture that every natural decidable problem can be solved in elementary time, but then you would be wrong. Consider the *extended regular expressions* defined by recursively combining (possibly empty) strings over some finite alphabet by concatenation, union (+), Kleene closure (\*), and negation. For example, the extended regular expression  $(0+1)^*00(0+1)^*$  represents the set of strings in  $\{0, 1\}^*$  that do *not* contain two 0s in a row. It is possible to determine algorithmically whether two extended regular expressions describe identical languages, by recursively converting each expression into an equivalent NFA, converting each NFA into a DFA, and then minimizing the DFA. Unfortunately, however, this equivalence problem cannot be decided in only elementary time!!

## Exercises

1. Consider the following problem, called **BoxDEPTH**: Given a set of  $n$  axis-aligned rectangles in the plane, how big is the largest subset of these rectangles that contain a common point?
  - (a) Describe a polynomial-time reduction from **BoxDEPTH** to **MAXCLIQUE**.
  - (b) Describe and analyze a polynomial-time algorithm for **BoxDEPTH**. [Hint:  $O(n^3)$  time should be easy, but  $O(n \log n)$  time is possible.]
  - (c) Why don't these two results imply that  $P=NP$ ?
2. (a) Describe a polynomial-time reduction from **PARTITION** to **SUBSETSUM**.
  - (b) Describe a polynomial-time reduction from **SUBSETSUM** to **PARTITION**.
3. (a) Describe and analyze an algorithm to solve **PARTITION** in time  $O(nM)$ , where  $n$  is the size of the input set and  $M$  is the sum of the absolute values of its elements.
  - (b) Why doesn't this algorithm imply that  $P=NP$ ?
4. A boolean formula is in *disjunctive normal form* (or *DNF*) if it consists of a *disjunction* (OR) or several *terms*, each of which is the conjunction (AND) of one or more literals. For example, the formula

$$(\bar{x} \wedge y \wedge \bar{z}) \vee (y \wedge z) \vee (x \wedge \bar{y} \wedge \bar{z})$$

is in disjunctive normal form. DNF-SAT asks, given a boolean formula in disjunctive normal form, whether that formula is satisfiable.

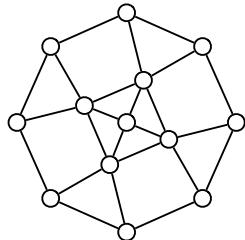
- (a) Describe a polynomial-time algorithm to solve DNF-SAT.
- (b) What is the error in the following argument that  $P=NP$ ?

Suppose we are given a boolean formula in conjunctive normal form with at most three literals per clause, and we want to know if it is satisfiable. We can use the distributive law to construct an equivalent formula in disjunctive normal form. For example,

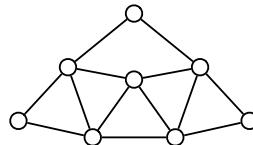
$$(x \vee y \vee \bar{z}) \wedge (\bar{x} \vee \bar{y}) \iff (x \wedge \bar{y}) \vee (y \wedge \bar{x}) \vee (\bar{z} \wedge \bar{x}) \vee (\bar{z} \wedge \bar{y})$$

Now we can use the algorithm from part (a) to determine, in polynomial time, whether the resulting DNF formula is satisfiable. We have just solved 3SAT in polynomial time. Since 3SAT is NP-hard, we must conclude that P=NP!

5. (a) Prove that PLANAR CIRCUITSAT is NP-complete.  
 (b) Prove that NOT ALL EQUAL 3SAT is NP-complete.  
 (c) Prove that the following variant of 3SAT is NP-complete: Given a boolean formula  $\Phi$  in conjunctive normal form where each clause contains at most 3 literals and each variable appears in at most 3 clauses, does  $\Phi$  have a satisfying assignment?
6. There's something special about the number 3.
  - (a) Describe and analyze a polynomial-time algorithm for 2PARTITION. Given a set  $S$  of  $2n$  positive integers, your algorithm will determine in polynomial time whether the elements of  $S$  can be split into  $n$  disjoint pairs whose sums are all equal.
  - (b) Describe and analyze a polynomial-time algorithm for 2COLOR. Given an undirected graph  $G$ , your algorithm will determine in polynomial time whether  $G$  has a proper coloring that uses only two colors.
  - (c) Describe and analyze a polynomial-time algorithm for 2SAT. Given a boolean formula  $\Phi$  in conjunctive normal form, with exactly two literals per clause, your algorithm will determine in polynomial time whether  $\Phi$  has a satisfying assignment.
7. There's nothing special about the number 3.
  - (a) The problem 12PARTITION is defined as follows: Given a set  $S$  of  $12n$  positive integers, determine whether the elements of  $S$  can be split into  $n$  subsets of 12 elements each whose sums are all equal. Prove that 12PARTITION is NP-hard. [Hint: Reduce from 3PARTITION. It may be easier to consider multisets first.]
  - (b) The problem 12COLOR is defined as follows: Given an undirected graph  $G$ , determine whether we can color each vertex with one of twelve colors, so that every edge touches two different colors. Prove that 12COLOR is NP-hard. [Hint: Reduce from 3COLOR.]
  - (c) The problem 12SAT is defined as follows: Given a boolean formula  $\Phi$  in conjunctive normal form, with exactly twelve literals per clause, determine whether  $\Phi$  has a satisfying assignment. Prove that 12SAT is NP-hard. [Hint: Reduce from 3SAT.]
8. (a) Using the gadget on the right below, prove that deciding whether a given planar graph is 3-colorable is NP-complete. [Hint: Show that the gadget can be 3-colored, and then replace any crossings in a planar embedding with the gadget appropriately.]  
 (b) Using part (a) and the middle gadget below, prove that deciding whether a planar graph with maximum degree 4 is 3-colorable is NP-complete. [Hint: Replace any vertex with degree greater than 4 with a collection of gadgets connected so that no degree is greater than four.]

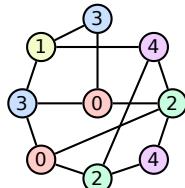


(a) Gadget for planar 3-colorability.



(b) Gadget for degree-4 planar 3-colorability.

9. Recall that a 5-coloring of a graph  $G$  is a function that assigns each vertex of  $G$  an ‘color’ from the set  $\{0, 1, 2, 3, 4\}$ , such that for any edge  $uv$ , vertices  $u$  and  $v$  are assigned different ‘colors’. A 5-coloring is *careful* if the colors assigned to adjacent vertices are not only distinct, but differ by more than 1 ( $\bmod 5$ ). Prove that deciding whether a given graph has a careful 5-coloring is NP-complete. [Hint: Reduce from the standard 5COLOR problem.]

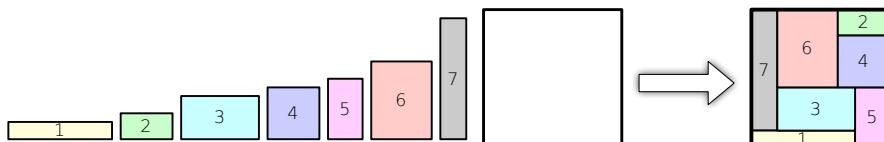


A careful 5-coloring.

10. Prove that the following problems are NP-complete.

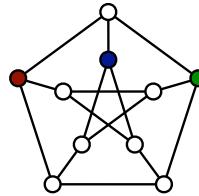
- (a) Given two undirected graphs  $G$  and  $H$ , is  $G$  isomorphic to a subgraph of  $H$ ?
- (b) Given an undirected graph  $G$ , does  $G$  have a spanning tree in which every node has degree at most 17?
- (c) Given an undirected graph  $G$ , does  $G$  have a spanning tree with at most 42 leaves?

11. The **RECTANGLETILING** problem is defined as follows: Given one large rectangle and several smaller rectangles, determine whether the smaller rectangles can be placed inside the large rectangle with no gaps or overlaps. Prove that **RECTANGLETILING** is NP-complete.

A positive instance of the **RECTANGLETILING** problem.

12. Let  $G = (V, E)$  be a graph. A *dominating set* in  $G$  is a subset  $S$  of the vertices such that every vertex in  $G$  is either in  $S$  or adjacent to a vertex in  $S$ . The **DOMINATINGSET** problem asks, given a graph  $G$  and an integer  $k$  as input, whether  $G$  contains a dominating set of size  $k$ . Either prove that this problem is NP-hard or describe a polynomial-time algorithm to solve it.

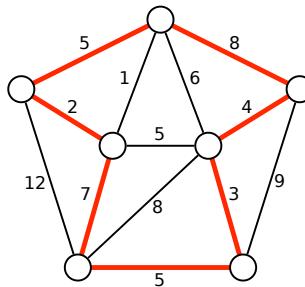
13. *Pebbling* is a solitaire game played on an undirected graph  $G$ , where each vertex has zero or more pebbles. A single *pebbling move* consists of removing two pebbles from a vertex  $v$  and adding one



A dominating set of size 3 in the Peterson graph.

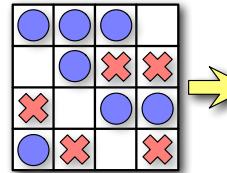
pebble to an arbitrary neighbor of  $v$ . (Obviously, the vertex  $v$  must have at least two pebbles before the move.) The PEBBLEDESTRUCTION problem asks, given a graph  $G = (V, E)$  and a pebble count  $p(v)$  for each vertex  $v$ , whether there is a sequence of pebbling moves that removes all but one pebble. Prove that PEBBLEDESTRUCTION is NP-complete.

14. (a) A *tonian path* in a graph  $G$  is a path that goes through at least half of the vertices of  $G$ . Show that determining whether a graph has a tonian path is NP-complete.  
 (b) A *tonian cycle* in a graph  $G$  is a cycle that goes through at least half of the vertices of  $G$ . Show that determining whether a graph has a tonian cycle is NP-complete. [Hint: Use part (a).]
  
15. For each problem below, either describe a polynomial-time algorithm or prove that the problem is NP-complete.
  - (a) A *double-Eulerian circuit* in an undirected graph  $G$  is a closed walk that traverses every edge in  $G$  exactly twice. Given a graph  $G$ , does  $G$  have a double-Eulerian circuit?
  - (b) A *double-Hamiltonian circuit* in an undirected graph  $G$  is a closed walk that visits every vertex in  $G$  exactly twice. Given a graph  $G$ , does  $G$  have a double-Hamiltonian circuit?
  
16. Let  $G$  be an undirected graph with weighted edges. A *heavy Hamiltonian cycle* is a cycle  $C$  that passes through each vertex of  $G$  exactly once, such that the total weight of the edges in  $C$  is at least half of the total weight of all edges in  $G$ . Prove that deciding whether a graph has a heavy Hamiltonian cycle is NP-complete.

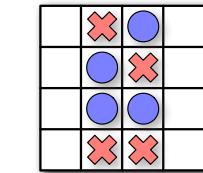
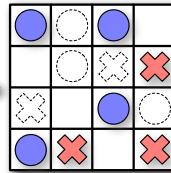


A heavy Hamiltonian cycle. The cycle has total weight 34; the graph has total weight 67.

17. Consider the following solitaire game. The puzzle consists of an  $n \times m$  grid of squares, where each square may be empty, occupied by a red stone, or occupied by a blue stone. The goal of the puzzle is to remove some of the given stones so that the remaining stones satisfy two conditions: (1) every row contains at least one stone, and (2) no column contains stones of both colors. For some initial configurations of stones, reaching this goal is impossible.



A solvable puzzle and one of its many solutions.



An unsolvable puzzle.

Prove that it is NP-hard to determine, given an initial configuration of red and blue stones, whether the puzzle can be solved.

18. A boolean formula in *exclusive-or conjunctive normal form* (XCNF) is a conjunction (AND) of several *clauses*, each of which is the *exclusive-or* of several literals; that is, a clause is true if and only if it contains an odd number of true literals. The XCNF-SAT problem asks whether a given XCNF formula is satisfiable. Either describe a polynomial-time algorithm for XCNF-SAT or prove that it is NP-hard.
19. Jeff tries to make his students happy. At the beginning of class, he passes out a questionnaire that lists a number of possible course policies in areas where he is flexible. Every student is asked to respond to each possible course policy with one of “strongly favor”, “mostly neutral”, or “strongly oppose”. Each student may respond with “strongly favor” or “strongly oppose” to at most five questions. Because Jeff’s students are very understanding, each student is happy if (but only if) he or she prevails in just one of his or her strong policy preferences. Either describe a polynomial-time algorithm for setting course policy to maximize the number of happy students, or show that the problem is NP-hard.
20. You’re in charge of choreographing a musical for your local community theater, and it’s time to figure out the final pose of the big show-stopping number at the end. (“Streetcar!”) You’ve decided that each of the  $n$  cast members in the show will be positioned in a big line when the song finishes, all with their arms extended and showing off their best spirit fingers.

The director has declared that during the final flourish, each cast member must either point both their arms up or point both their arms down; it’s your job to figure out who points up and who points down. Moreover, in a fit of unchecked power, the director has also given you a list of arrangements that will upset his delicate artistic temperament. Each forbidden arrangement is a subset of the cast members paired with arm positions; for example: “Marge may not point her arms up while Ned, Apu, and Smithers point their arms down.”

Prove that finding an acceptable arrangement of arm positions is NP-hard.

21. The next time you are at a party, one of the guests will suggest everyone play a round of Three-Way Mumbledypeg, a game of skill and dexterity that requires three teams and a knife. The official Rules of Three-Way Mumbledypeg (fixed during the Holy Roman Three-Way Mumbledypeg Council in 1625) require that (1) each team *must* have at least one person, (2) any two people on the same team *must* know each other, and (3) everyone watching the game *must* be on one of the three teams. Of course, it will be a really *fun* party; nobody will want to leave. There will be several pairs of people at the party who don’t know each other. The host of the party, having

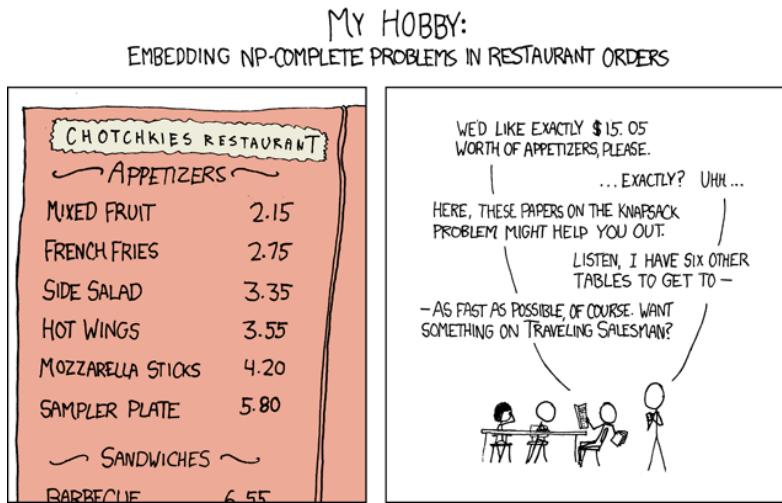
heard thrilling tales of your prowess in all things algorithmic, will hand you a list of which pairs of party-goers know each other and ask you to choose the teams, while he sharpens the knife.

Either describe and analyze a polynomial time algorithm to determine whether the party-goers can be split into three legal Three-Way Mumbledypeg teams, or prove that the problem is NP-hard.

22. The party you are attending is going great, but now it's time to line up for *The Algorithm March* (アルゴリズムこうしん)! This dance was originally developed by the Japanese comedy duo Itsumo Kokokara (いつもここから) for the children's television show PythagoraSwitch (ピタゴラスイッチ). The Algorithm March is performed by a line of people; each person in line starts a specific sequence of movements one measure later than the person directly in front of them. Thus, the march is the dance equivalent of a musical round or canon, like "Row Row Row Your Boat".

Proper etiquette dictates that each marcher must know the person directly in front of them in line, lest a minor mistake during lead to horrible embarrassment between strangers. Suppose you are given a complete list of which people at your party know each other. **Prove** that it is NP-hard to determine the largest number of party-goers that can participate in the Algorithm March. You may assume there are no ninjas at your party.

23. (a) Suppose you are given a magic black box that can determine **in polynomial time**, given an arbitrary weighted graph  $G$ , the length of the shortest Hamiltonian cycle in  $G$ . Describe and analyze a **polynomial-time** algorithm that computes, given an arbitrary weighted graph  $G$ , the shortest Hamiltonian cycle in  $G$ , using this magic black box as a subroutine.
- (b) Suppose you are given a magic black box that can determine **in polynomial time**, given an arbitrary graph  $G$ , the number of vertices in the largest complete subgraph of  $G$ . Describe and analyze a **polynomial-time** algorithm that computes, given an arbitrary graph  $G$ , a complete subgraph of  $G$  of maximum size, using this magic black box as a subroutine.
- (c) Suppose you are given a magic black box that can determine **in polynomial time**, given an arbitrary weighted graph  $G$ , whether  $G$  is 3-colorable. Describe and analyze a **polynomial-time** algorithm that either computes a proper 3-coloring of a given graph or correctly reports that no such coloring exists, using the magic black box as a subroutine. [*Hint: The input to the magic black box is a graph. Just a graph. Vertices and edges. Nothing else.*]
- (d) Suppose you are given a magic black box that can determine **in polynomial time**, given an arbitrary boolean formula  $\Phi$ , whether  $\Phi$  is satisfiable. Describe and analyze a **polynomial-time** algorithm that either computes a satisfying assignment for a given boolean formula or correctly reports that no such assignment exists, using the magic black box as a subroutine.
- \*(e) Suppose you are given a magic black box that can determine **in polynomial time**, given an initial Tetris configuration and a finite sequence of Tetris pieces, whether a perfect player can play every piece. (This problem is NP-hard.) Describe and analyze a **polynomial-time** algorithm that computes the shortest Hamiltonian cycle in a given weighted graph in polynomial time, using this magic black box as a subroutine. [*Hint: Use part (a). You do not need to know anything about Tetris to solve this problem.*]



[General solutions give you a 50% tip.]

— Randall Munroe, [xkcd \(<http://xkcd.com/287/>\)](http://xkcd.com/287/)  
Reproduced under a Creative Commons Attribution-NonCommercial 2.5 License

*Le mieux est l'ennemi du bien. [The best is the enemy of the good.]*

— Voltaire, *La Béguerie* (1772)

*Who shall forbid a wise skepticism, seeing that there is no practical question on which any thing more than an approximate solution can be had?*

— Ralph Waldo Emerson, *Representative Men* (1850)

*Now, distrust of corporations threatens our still-tentative economic recovery; it turns out greed is bad, after all.*

— Paul Krugman, "Greed is Bad", *The New York Times*, June 4, 2002.

## \*30 Approximation Algorithms

### 30.1 Load Balancing

On the future smash hit reality-TV game show *Grunt Work*, scheduled to air Thursday nights at 3am (2am Central) on ESPN $\pi$ , the contestants are given a series of utterly pointless tasks to perform. Each task has a predetermined time limit; for example, "Sharpen this pencil for 17 seconds", or "Pour pig's blood on your head and sing The Star-Spangled Banner for two minutes", or "Listen to this 75-minute algorithms lecture". The directors of the show want you to assign each task to one of the contestants, so that the last task is completed as early as possible. When your predecessor correctly informed the directors that their problem is NP-hard, he was immediately fired. "Time is money!" they screamed at him. "We don't need perfection. Wake up, dude, this is television!"

Less facetiously, suppose we have a set of  $n$  jobs, which we want to assign to  $m$  machines. We are given an array  $T[1..n]$  of non-negative numbers, where  $T[j]$  is the running time of job  $j$ . We can describe an *assignment* by an array  $A[1..n]$ , where  $A[j] = i$  means that job  $j$  is assigned to machine  $i$ . The *makespan* of an assignment is the maximum time that any machine is busy:

$$\text{makespan}(A) = \max_i \sum_{A[j]=i} T[j]$$

The *load balancing* problem is to compute the assignment with the smallest possible makespan.

It's not hard to prove that the load balancing problem is NP-hard by reduction from PARTITION: The array  $T[1..n]$  can be evenly partitioned if and only if there is an assignment to two machines with makespan exactly  $\sum_i T[i]/2$ . A slightly more complicated reduction from 3PARTITION implies that the load balancing problem is *strongly* NP-hard. If we really need the optimal solution, there is a dynamic programming algorithm that runs in time  $O(nM^m)$ , where  $M$  is the minimum makespan, but that's just horrible.

There is a fairly natural and efficient greedy heuristic for load balancing: consider the jobs one at a time, and assign each job to the machine  $i$  with the earliest finishing time  $Total[i]$ .

```

GREEDYLOADBALANCE( $T[1..n], m$ ):
  for  $i \leftarrow 1$  to  $m$ 
     $Total[i] \leftarrow 0$ 
  for  $j \leftarrow 1$  to  $n$ 
     $mini \leftarrow \arg \min_i Total[i]$ 
     $A[j] \leftarrow mini$ 
     $Total[mini] \leftarrow Total[mini] + T[j]$ 
  return  $A[1..m]$ 

```

**Theorem 1.** *The makespan of the assignment computed by GREEDYLOADBALANCE is at most twice the makespan of the optimal assignment.*

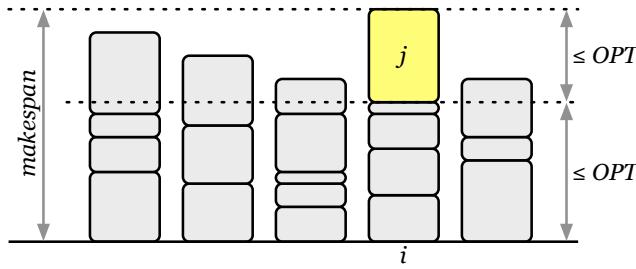
**Proof:** Fix an arbitrary input, and let  $OPT$  denote the makespan of its optimal assignment. The approximation bound follows from two trivial observations. First, the makespan of any assignment (and therefore of the optimal assignment) is at least the duration of the longest job. Second, the makespan of any assignment is at least the total duration of all the jobs divided by the number of machines.

$$OPT \geq \max_j T[j] \quad \text{and} \quad OPT \geq \frac{1}{m} \sum_{j=1}^n T[j]$$

Now consider the assignment computed by GREEDYLOADBALANCE. Suppose machine  $i$  has the largest total running time, and let  $j$  be the last job assigned to machine  $i$ . Our first trivial observation implies that  $T[j] \leq OPT$ . To finish the proof, we must show that  $Total[i] - T[j] \leq OPT$ . Job  $j$  was assigned to machine  $i$  because it had the smallest finishing time, so  $Total[i] - T[j] \leq Total[k]$  for all  $k$ . (Some values  $Total[k]$  may have increased since job  $j$  was assigned, but that only helps us.) In particular,  $Total[i] - T[j]$  is less than or equal to the average finishing time over all machines. Thus,

$$Total[i] - T[j] \leq \frac{1}{m} \sum_{i=1}^m Total[i] = \frac{1}{m} \sum_{j=1}^n T[j] \leq OPT$$

by our second trivial observation. We conclude that the makespan  $Total[i]$  is at most  $2 \cdot OPT$ .  $\square$



Proof that GREEDYLOADBALANCE is a 2-approximation algorithm

GREEDYLOADBALANCE is an *online* algorithm: It assigns jobs to machines in the order that the jobs appear in the input array. Online approximation algorithms are useful in settings where inputs arrive in a stream of unknown length—for example, real jobs arriving at a real scheduling algorithm. In this online setting, it may be *impossible* to compute an optimum solution, even in cases where the offline problem (where all inputs are known in advance) can be solved in polynomial time. The study of online algorithms could easily fill an entire one-semester course (alas, not this one).

In our original offline setting, we can improve the approximation factor by sorting the jobs before piping them through the greedy algorithm.

<b>SORTEDGREEDYLOADBALANCE(<math>T[1..n], m</math>):</b> sort $T$ in decreasing order return GREEDYLOADBALANCE( $T, m$ )
--

**Theorem 2.** *The makespan of the assignment computed by SORTEDGREEDYLOADBALANCE is at most  $3/2$  times the makespan of the optimal assignment.*

**Proof:** Let  $i$  be the busiest machine in the schedule computed by SORTEDGREEDYLOADBALANCE. If only one job is assigned to machine  $i$ , then the greedy schedule is actually optimal, and the theorem is trivially true. Otherwise, let  $j$  be the last job assigned to machine  $i$ . Since each of the first  $m$  jobs is assigned to a unique machine, we must have  $j \geq m + 1$ . As in the previous proof, we know that  $\text{Total}[i] - T[j] \leq OPT$ .

In any schedule, at least two of the first  $m + 1$  jobs, say jobs  $k$  and  $\ell$ , must be assigned to the same machine. Thus,  $T[k] + T[\ell] \leq OPT$ . Since  $\max\{k, \ell\} \leq m + 1 \leq j$ , and the jobs are sorted in decreasing order by duration, we have

$$T[j] \leq T[m + 1] \leq T[\max\{k, \ell\}] = \min\{T[k], T[\ell]\} \leq OPT/2.$$

We conclude that the makespan  $\text{Total}[i]$  is at most  $3 \cdot OPT/2$ , as claimed.  $\square$

## 30.2 Generalities

Consider an arbitrary optimization problem. Let  $OPT(X)$  denote the value of the optimal solution for a given input  $X$ , and let  $A(X)$  denote the value of the solution computed by algorithm  $A$  given the same input  $X$ . We say that  $A$  is an  **$\alpha(n)$ -approximation algorithm** if and only if

$$\frac{OPT(X)}{A(X)} \leq \alpha(n) \quad \text{and} \quad \frac{A(X)}{OPT(X)} \leq \alpha(n)$$

for all inputs  $X$  of size  $n$ . The function  $\alpha(n)$  is called the **approximation factor** for algorithm  $A$ . For any given algorithm, only one of these two inequalities will be important. For maximization problems, where we want to compute a solution whose cost is as small as possible, the first inequality is trivial. For maximization problems, where we want a solution whose value is as large as possible, the second inequality is trivial. A 1-approximation algorithm always returns the exact optimal solution.

Especially for problems where exact optimization is NP-hard, we have little hope of completely characterizing the optimal solution. The secret to proving that an algorithm satisfies some approximation ratio is to find a useful function of the input that provides both lower bounds on the cost of the optimal solution and upper bounds on the cost of the approximate solution. For example, if  $OPT(X) \geq f(X)/2$  and  $A(X) \leq 5f(X)$  for any function  $f$ , then  $A$  is a 10-approximation algorithm. Finding the right intermediate function can be a delicate balancing act.

## 30.3 Greedy Vertex Cover

Recall that the *vertex cover* problem asks, given a graph  $G$ , for the smallest set of vertices of  $G$  that cover every edge. This is one of the first NP-hard problems introduced in the first week of class. There is a natural and efficient greedy heuristic<sup>1</sup> for computing a small vertex cover: mark the vertex with the largest degree, remove all the edges incident to that vertex, and recurse.

```

GREEDYVERTEXCOVER( $G$ ):
   $C \leftarrow \emptyset$ 
  while  $G$  has at least one edge
     $v \leftarrow$  vertex in  $G$  with maximum degree
     $G \leftarrow G \setminus v$ 
     $C \leftarrow C \cup v$ 
  return  $C$ 

```

Obviously this algorithm doesn't compute the optimal vertex cover—that would imply P=NP!—but it does compute a reasonably close approximation.

<sup>1</sup>A *heuristic* is an algorithm that doesn't work.

**Theorem 3.** GREEDYVERTEXCOVER is an  $O(\log n)$ -approximation algorithm.

**Proof:** For all  $i$ , let  $G_i$  denote the graph  $G$  after  $i$  iterations of the main loop, and let  $d_i$  denote the maximum degree of any node in  $G_{i-1}$ . We can define these variables more directly by adding a few extra lines to our algorithm:

```

GREEDYVERTEXCOVER( $G$ ):
   $C \leftarrow \emptyset$ 
   $G_0 \leftarrow G$ 
   $i \leftarrow 0$ 
  while  $G_i$  has at least one edge
     $i \leftarrow i + 1$ 
     $v_i \leftarrow$  vertex in  $G_{i-1}$  with maximum degree
     $d_i \leftarrow \deg_{G_{i-1}}(v_i)$ 
     $G_i \leftarrow G_{i-1} \setminus v_i$ 
     $C \leftarrow C \cup v_i$ 
  return  $C$ 
```

Let  $|G_{i-1}|$  denote the number of edges in the graph  $G_{i-1}$ . Let  $C^*$  denote the optimal vertex cover of  $G$ , which consists of  $OPT$  vertices. Since  $C^*$  is also a vertex cover for  $G_{i-1}$ , we have

$$\sum_{v \in C^*} \deg_{G_{i-1}}(v) \geq |G_{i-1}|.$$

In other words, the *average* degree in  $G_i$  of any node in  $C^*$  is at least  $|G_{i-1}|/OPT$ . It follows that  $G_{i-1}$  has at least one node with degree at least  $|G_{i-1}|/OPT$ . Since  $d_i$  is the maximum degree of any node in  $G_{i-1}$ , we have

$$d_i \geq \frac{|G_{i-1}|}{OPT}$$

Moreover, for any  $j \geq i - 1$ , the subgraph  $G_j$  has no more edges than  $G_{i-1}$ , so  $d_i \geq |G_j|/OPT$ . This observation implies that

$$\sum_{i=1}^{OPT} d_i \geq \sum_{i=1}^{OPT} \frac{|G_{i-1}|}{OPT} \geq \sum_{i=1}^{OPT} \frac{|G_{OPT}|}{OPT} = |G_{OPT}| = |G| - \sum_{i=1}^{OPT} d_i.$$

In other words, the first  $OPT$  iterations of GREEDYVERTEXCOVER remove at least half the edges of  $G$ . Thus, after at most  $OPT \lg |G| \leq 2 OPT \lg n$  iterations, all the edges of  $G$  have been removed, and the algorithm terminates. We conclude that GREEDYVERTEXCOVER computes a vertex cover of size  $O(OPT \log n)$ .  $\square$

So far we've only proved an *upper bound* on the approximation factor of GREEDYVERTEXCOVER; perhaps a more careful analysis would imply that the approximation factor is only  $O(\log \log n)$ , or even  $O(1)$ . Alas, no such improvement is possible. For any integer  $n$ , a simple recursive construction gives us an  $n$ -vertex graph for which the greedy algorithm returns a vertex cover of size  $\Omega(OPT \cdot \log n)$ . Details are left as an exercise for the reader.

### 30.4 Set Cover and Hitting Set

The greedy algorithm for vertex cover can be applied almost immediately to two more general problems: *set cover* and *hitting set*. The input for both of these problems is a *set system*  $(X, \mathcal{F})$ , where  $X$  is a finite *ground set*, and  $\mathcal{F}$  is a family of subsets of  $X$ .<sup>2</sup> A *set cover* of a set system  $(X, \mathcal{F})$  is a subfamily of sets in

<sup>2</sup>A matroid (see the lecture note on greedy algorithms) is a special type of set system.

$\mathcal{F}$  whose union is the entire ground set  $X$ . A *hitting set* for  $(X, \mathcal{F})$  is a subset of the ground set  $X$  that intersects every set in  $\mathcal{F}$ .

An undirected graph can be cast as a set system in two different ways. In one formulation, the ground set  $X$  contains the vertices, and each edge defines a set of two vertices in  $\mathcal{F}$ . In this formulation, a vertex cover is a hitting set. In the other formulation, the *edges* are the ground set, the *vertices* define the family of subsets, and a vertex cover is a set cover.

Here are the natural greedy algorithms for finding a small set cover and finding a small hitting set. GREEDYSETCOVER finds a set cover whose size is at most  $O(\log |\mathcal{F}|)$  times the size of smallest set cover. GREEDYHITTINGSET finds a hitting set whose size is at most  $O(\log |X|)$  times the size of the smallest hitting set.

```
GREEDYSETCOVER( $X, \mathcal{F}$ ):
 $\mathcal{C} \leftarrow \emptyset$ 
while  $X \neq \emptyset$ 
     $S \leftarrow \arg \max_{S \in \mathcal{F}} |S \cap X|$ 
     $X \leftarrow X \setminus S$ 
     $\mathcal{C} \leftarrow \mathcal{C} \cup \{S\}$ 
return  $\mathcal{C}$ 
```

```
GREEDYHITTINGSET( $X, \mathcal{F}$ ):
 $H \leftarrow \emptyset$ 
while  $\mathcal{F} \neq \emptyset$ 
     $x \leftarrow \arg \max_{x \in X} |\{S \in \mathcal{F} \mid x \in S\}|$ 
     $\mathcal{F} \leftarrow \mathcal{F} \setminus \{S \in \mathcal{F} \mid x \in S\}$ 
     $H \leftarrow H \cup \{x\}$ 
return  $H$ 
```

The similarity between these two algorithms is no coincidence. For any set system  $(X, \mathcal{F})$ , there is a *dual* set system  $(\mathcal{F}, X^*)$  defined as follows. For any element  $x \in X$  in the ground set, let  $x^*$  denote the subfamily of sets in  $\mathcal{F}$  that contain  $x$ :

$$x^* = \{S \in \mathcal{F} \mid x \in S\}.$$

Finally, let  $X^*$  denote the collection of all subsets of the form  $x^*$ :

$$X^* = \{x^* \mid x \in X\}.$$

As an example, suppose  $X$  is the set of letters of alphabet and  $\mathcal{F}$  is the set of last names of student taking CS 573 this semester. Then  $X^*$  has 26 elements, each containing the subset of CS 573 students whose last name contains a particular letter of the alphabet. For example,  $m^*$  is the set of students whose last names contain the letter  $m$ .

There is a direct one-to-one correspondence between the ground set  $X$  and the dual set family  $X^*$ . It is a tedious but instructive exercise to prove that the dual of the dual of any set system is isomorphic to the original set system— $(X^*, \mathcal{F}^*)$  is essentially the same as  $(X, \mathcal{F})$ . It is also easy to prove that a set cover for any set system  $(X, \mathcal{F})$  is also a hitting set for the dual set system  $(\mathcal{F}, X^*)$ , and therefore a hitting set for any set system  $(X, \mathcal{F})$  is isomorphic to a set cover for the dual set system  $(\mathcal{F}, X^*)$ .

### 30.5 Vertex Cover, Again

The greedy approach doesn't always lead to the best approximation algorithms. Consider the following alternate heuristic for vertex cover:

```
DUMBVERTEXCOVER( $G$ ):
 $C \leftarrow \emptyset$ 
while  $G$  has at least one edge
     $(u, v) \leftarrow$  any edge in  $G$ 
     $G \leftarrow G \setminus \{u, v\}$ 
     $C \leftarrow C \cup \{u, v\}$ 
return  $C$ 
```

The minimum vertex cover—in fact, *every* vertex cover—contains at least one of the two vertices  $u$  and  $v$  chosen inside the while loop. It follows immediately that DUMBVERTEXCOVER is a 2-approximation algorithm!

The same idea can be extended to approximate the minimum hitting set for any set system  $(X, \mathcal{F})$ , where the approximation factor is the size of the largest set in  $\mathcal{F}$ .

### 30.6 Traveling Salesman: The Bad News

The *traveling salesman problem*<sup>3</sup> problem asks for the shortest Hamiltonian cycle in a weighted undirected graph. To keep the problem simple, we can assume without loss of generality that the underlying graph is always the complete graph  $K_n$  for some integer  $n$ ; thus, the input to the traveling salesman problem is just a list of the  $\binom{n}{2}$  edge lengths.

Not surprisingly, given its similarity to the Hamiltonian cycle problem, it's quite easy to prove that the traveling salesman problem is NP-hard. Let  $G$  be an arbitrary undirected graph with  $n$  vertices. We can construct a length function for  $K_n$  as follows:

$$\ell(e) = \begin{cases} 1 & \text{if } e \text{ is an edge in } G, \\ 2 & \text{otherwise.} \end{cases}$$

Now it should be obvious that if  $G$  has a Hamiltonian cycle, then there is a Hamiltonian cycle in  $K_n$  whose length is exactly  $n$ ; otherwise every Hamiltonian cycle in  $K_n$  has length at least  $n + 1$ . We can clearly compute the lengths in polynomial time, so we have a polynomial time reduction from Hamiltonian cycle to traveling salesman. Thus, the traveling salesman problem is NP-hard, even if all the edge lengths are 1 and 2.

There's nothing special about the values 1 and 2 in this reduction; we can replace them with any values we like. By choosing values that are sufficiently far apart, we can show that even approximating the shortest traveling salesman tour is NP-hard. For example, suppose we set the length of the ‘absent’ edges to  $n + 1$  instead of 2. Then the shortest traveling salesman tour in the resulting weighted graph either has length exactly  $n$  (if  $G$  has a Hamiltonian cycle) or has length at least  $2n$  (if  $G$  does not have a Hamiltonian cycle). Thus, if we could approximate the shortest traveling salesman tour within a factor of 2 in polynomial time, we would have a polynomial-time algorithm for the Hamiltonian cycle problem.

Pushing this idea to its limits us the following negative result.

**Theorem 4.** *For any function  $f(n)$  that can be computed in time polynomial in  $n$ , there is no polynomial-time  $f(n)$ -approximation algorithm for the traveling salesman problem on general weighted graphs, unless  $P=NP$ .*

### 30.7 Traveling Salesman: The Good News

Even though the general traveling salesman problem can't be approximated, a common special case can be approximated fairly easily. The special case requires the edge lengths to satisfy the so-called *triangle inequality*:

$$\ell(u, w) \leq \ell(u, v) + \ell(v, w) \quad \text{for any vertices } u, v, w.$$

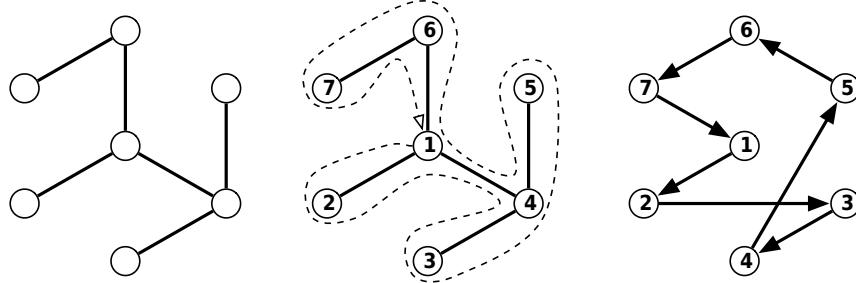
This inequality is satisfied for *geometric graphs*, where the vertices are points in the plane (or some higher-dimensional space), edges are straight line segments, and lengths are measured in the usual

---

<sup>3</sup>This is sometimes bowdlerized into the traveling salesperson problem. That's just silly. Who ever heard of a traveling salesperson sleeping with the farmer's child?

Euclidean metric. Notice that the length functions we used above to show that the general TSP is hard to approximate do not (always) satisfy the triangle inequality.

With the triangle inequality in place, we can quickly compute a 2-approximation for the traveling salesman tour as follows. First, we compute the minimum spanning tree  $T$  of the weighted input graph; this can be done in  $O(n^2 \log n)$  time (where  $n$  is the number of vertices of the graph) using any of several classical algorithms. Second, we perform a depth-first traversal of  $T$ , numbering the vertices in the order that we first encounter them. Because  $T$  is a spanning tree, every vertex is numbered. Finally, we return the cycle obtained by visiting the vertices according to this numbering.



A minimum spanning tree  $T$ , a depth-first traversal of  $T$ , and the resulting approximate traveling salesman tour.

**Theorem 5.** *A depth-first ordering of the minimum spanning tree gives a 2-approximation of the shortest traveling salesman tour.*

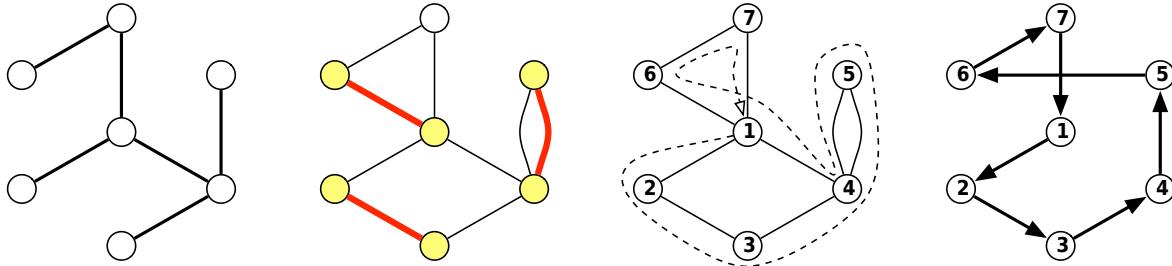
**Proof:** Let  $OPT$  denote the cost of the optimal TSP tour, let  $MST$  denote the total length of the minimum spanning tree, and let  $A$  be the length of the tour computed by our approximation algorithm. Consider the ‘tour’ obtained by walking through the minimum spanning tree in depth-first order. Since this tour traverses every edge in the tree exactly twice, its length is  $2 \cdot MST$ . The final tour can be obtained from this one by removing duplicate vertices, moving directly from each node to the next *unvisited* node.; the triangle inequality implies that taking these shortcuts cannot make the tour longer. Thus,  $A \leq 2 \cdot MST$ . On the other hand, if we remove any edge from the optimal tour, we obtain a spanning tree (in fact a spanning path) of the graph; thus,  $MST \geq OPT$ . We conclude that  $A \leq 2 \cdot OPT$ ; our algorithm computes a 2-approximation of the optimal tour.  $\square$

We can improve this approximation factor using the following algorithm discovered by Nicos Christofides in 1976. As in the previous algorithm, we start by constructing the minimum spanning tree  $T$ . Then let  $O$  be the set of vertices with odd degree in  $T$ ; it is an easy exercise (hint, hint) to show that the number of vertices in  $O$  is even.

In the next stage of the algorithm, we compute a *minimum-cost perfect matching*  $M$  of these odd-degree vertices. A perfect matching is a collection of edges, where each edge has both endpoints in  $O$  and each vertex in  $O$  is adjacent to exactly one edge; we want the perfect matching of minimum total length. Later in the semester, we will see an algorithm to compute  $M$  in polynomial time.

Now consider the multigraph  $T \cup M$ ; any edge in both  $T$  and  $M$  appears twice in this multigraph. This graph is connected, and every vertex has even degree. Thus, it contains an *Eulerian circuit*: a closed walk that uses every edge exactly once. We can compute such a walk in  $O(n)$  time with a simple modification of depth-first search. To obtain the final approximate TSP tour, we number the vertices in the order they first appear in some Eulerian circuit of  $T \cup M$ , and return the cycle obtained by visiting the vertices according to that numbering.

**Theorem 6.** *Given a weighted graph that obeys the triangle inequality, the Christofides heuristic computes a  $(3/2)$ -approximation of the shortest traveling salesman tour.*



A minimum spanning tree  $T$ , a minimum-cost perfect matching  $M$  of the odd vertices in  $T$ , an Eulerian circuit of  $T \cup M$ , and the resulting approximate traveling salesman tour.

**Proof:** Let  $A$  denote the length of the tour computed by the Christofides heuristic; let  $OPT$  denote the length of the optimal tour; let  $MST$  denote the total length of the minimum spanning tree; let  $MOM$  denote the total length of the minimum odd-vertex matching.

The graph  $T \cup M$ , and therefore any Euler tour of  $T \cup M$ , has total length  $MST + MOM$ . By the triangle inequality, taking a shortcut past a previously visited vertex can only shorten the tour. Thus,  $A \leq MST + MOM$ .

By the triangle inequality, the optimal tour of the odd-degree vertices of  $T$  cannot be longer than  $OPT$ . Any cycle passing through of the odd vertices can be partitioned into two perfect matchings, by alternately coloring the edges of the cycle red and green. One of these two matchings has length at most  $OPT/2$ . On the other hand, both matchings have length at least  $MOM$ . Thus,  $MOM \leq OPT/2$ .

Finally, recall our earlier observation that  $MST \leq OPT$ .

Putting these three inequalities together, we conclude that  $A \leq 3 \cdot OPT/2$ , as claimed.  $\square$

### 30.8 $k$ -center Clustering

The  $k$ -center clustering problem is defined as follows. We are given a set  $P = \{p_1, p_2, \dots, p_n\}$  of  $n$  points in the plane<sup>4</sup> and an integer  $k$ . Our goal to find a collection of  $k$  circles that collectively enclose all the input points, such that the radius of the largest circle is as large as possible. More formally, we want to compute a set  $C = \{c_1, c_2, \dots, c_k\}$  of  $k$  center points, such that the following cost function is minimized:

$$\text{cost}(C) := \max_i \min_j |p_i c_j|.$$

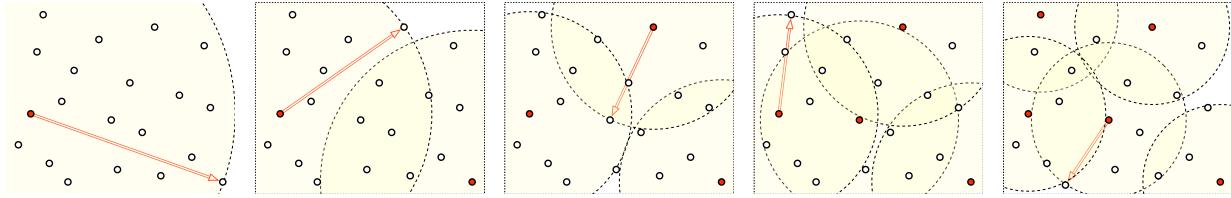
Here,  $|p_i c_j|$  denotes the Euclidean distance between input point  $p_i$  and center point  $c_j$ . Intuitively, each input point is assigned to its closest center point; the points assigned to a given center  $c_j$  comprise a *cluster*. The distance from  $c_j$  to the farthest point in its cluster is the *radius* of that cluster; the cluster is contained in a circle of this radius centered at  $c_j$ . The  $k$ -center clustering cost  $\text{cost}(C)$  is precisely the maximum cluster radius.

This problem turns out to be NP-hard, even to approximate within a factor of roughly 1.8. However, there is a natural greedy strategy, first analyzed in 1985 by Teofilo Gonzalez<sup>5</sup>, that is guaranteed to produce a clustering whose cost is at most twice optimal. Choose the  $k$  center points one at a time, starting with an arbitrary input point as the first center. In each iteration, choose the input point that is farthest from any earlier center point to be the next center point.

In the pseudocode below,  $d_i$  denotes the current distance from point  $p_i$  to its nearest center, and  $r_j$  denotes the maximum of all  $d_i$  (or in other words, the cluster radius) after the first  $j$  centers have

<sup>4</sup>The  $k$ -center problem can be defined over *any* metric space, and the approximation analysis in this section holds in any metric space as well. The analysis in the *next* section, however, does require that the points come from the Euclidean plane.

<sup>5</sup>Teofilo F. Gonzalez. Clustering to minimize the maximum inter-cluster distance. *Theoretical Computer Science* 38:293–306, 1985.

The first five iterations of Gonzalez's  $k$ -center clustering algorithm.

been chosen. The algorithm includes an extra iteration to compute the final clustering radius  $r_k$  (and the next center  $c_{k+1}$ ).

```
GONZALEZKCENTER( $P, k$ ):
  for  $i \leftarrow 1$  to  $n$ 
     $d_i \leftarrow \infty$ 
   $c_1 \leftarrow p_1$ 
  for  $j \leftarrow 1$  to  $k$ 
     $r_j \leftarrow 0$ 
    for  $i \leftarrow 1$  to  $n$ 
       $d_i \leftarrow \min\{d_i, |p_i c_j|\}$ 
      if  $r_j < d_i$ 
         $r_j \leftarrow d_i; c_{j+1} \leftarrow p_i$ 
  return  $\{c_1, c_2, \dots, c_k\}$ 
```

GONZALEZKCENTER clearly runs in  $O(nk)$  time. Using more advanced data structures, Tomas Feder and Daniel Greene<sup>6</sup> described an algorithm to compute exactly the same clustering in only  $O(n \log k)$  time.

**Theorem 7.** GONZALEZKCENTER computes a 2-approximation to the optimal  $k$ -center clustering.

**Proof:** Let  $OPT$  denote the optimal  $k$ -center clustering radius for  $P$ . For any index  $i$ , let  $c_i$  and  $r_i$  denote the  $i$ th center point and  $i$ th clustering radius computed by GONZALEZKCENTER.

By construction, each center point  $c_j$  has distance at least  $r_{j-1}$  from any center point  $c_i$  with  $i < j$ . Moreover, for any  $i < j$ , we have  $r_i \geq r_j$ . Thus,  $|c_i c_j| \geq r_k$  for all indices  $i$  and  $j$ .

On the other hand, at least one cluster in the optimal clustering contains at least two of the points  $c_1, c_2, \dots, c_{k+1}$ . Thus, by the triangle inequality, we must have  $|c_i c_j| \leq 2 \cdot OPT$  for some indices  $i$  and  $j$ . We conclude that  $r_k \leq 2 \cdot OPT$ , as claimed.  $\square$

### \*30.9 Approximation Schemes

With just a little more work, we can compute an arbitrarily close approximation of the optimal  $k$ -clustering, using a so-called *approximation scheme*. An approximation scheme accepts both an instance of the problem and a value  $\epsilon > 0$  as input, and it computes a  $(1 + \epsilon)$ -approximation of the optimal output for that instance. As I mentioned earlier, computing even a 1.8-approximation is NP-hard, so we cannot expect our approximation scheme to run in polynomial time; nevertheless, at least for small values of  $k$ , the approximation scheme will be considerably more efficient than any exact algorithm.

Our approximation scheme works in three phases:

<sup>6</sup>Tomas Feder\* and Daniel H. Greene. Optimal algorithms for approximate clustering. *Proc. 20th STOC*, 1988. Unlike Gonzalez's algorithm, Feder and Greene's faster algorithm does not work over arbitrary metric spaces; it requires that the input points come from some  $\mathbb{R}^d$  and that distances are measured in some  $L_p$  metric. The time analysis also assumes that the distance between any two points can be computed in  $O(1)$  time.

1. Compute a 2-approximate clustering of the input set  $P$  using GONZALEZKCENTER. Let  $r$  be the cost of this clustering.
2. Create a regular grid of squares of width  $\delta = \varepsilon r / 2\sqrt{2}$ . Let  $Q$  be a subset of  $P$  containing one point from each non-empty cell of this grid.
3. Compute an *optimal* set of  $k$  centers for  $Q$ . Return these  $k$  centers as the approximate  $k$ -center clustering for  $P$ .

The first phase requires  $O(nk)$  time. By our earlier analysis, we have  $r^* \leq r \leq 2r^*$ , where  $r^*$  is the optimal  $k$ -center clustering cost for  $P$ .

The second phase can be implemented in  $O(n)$  time using a hash table, or in  $O(n \log n)$  time by standard sorting, by associating approximate coordinates  $(\lfloor x/\delta \rfloor, \lfloor y/\delta \rfloor)$  to each point  $(x, y) \in P$  and removing duplicates. The key observation is that the resulting point set  $Q$  is significantly smaller than  $P$ . We know  $P$  can be covered by  $k$  balls of radius  $r^*$ , each of which touches  $O(r^*/\delta^2) = O(1/\varepsilon^2)$  grid cells. It follows that  $|Q| = O(k/\varepsilon^2)$ .

Let  $T(n, k)$  be the running time of an *exact*  $k$ -center clustering algorithm, given  $n$  points as input. If this were a computational geometry class, we might see a “brute force” algorithm that runs in time  $T(n, k) = O(n^{k+2})$ ; the fastest algorithm currently known<sup>7</sup> runs in time  $T(n, k) = n^{O(\sqrt{k})}$ . If we use this algorithm, our third phase requires  $(k/\varepsilon^2)^{O(\sqrt{k})}$  time.

It remains to show that the optimal clustering for  $Q$  implies a  $(1 + \varepsilon)$ -approximation of the optimal clustering for  $P$ . Suppose the optimal clustering of  $Q$  consists of  $k$  balls  $B_1, B_2, \dots, B_k$ , each of radius  $\tilde{r}$ . Clearly  $\tilde{r} \leq r^*$ , since any set of  $k$  balls that cover  $P$  also cover any subset of  $P$ . Each point in  $P \setminus Q$  shares a grid cell with some point in  $Q$ , and therefore is within distance  $\delta\sqrt{2}$  of some point in  $Q$ . Thus, if we increase the radius of each ball  $B_i$  by  $\delta\sqrt{2}$ , the expanded balls must contain every point in  $P$ . We conclude that the optimal centers for  $Q$  gives us a  $k$ -center clustering for  $P$  of cost at most  $r^* + \delta\sqrt{2} \leq r^* + \varepsilon r^* = (1 + \varepsilon)r^*$ .

The total running time of the approximation scheme is  $O(nk + (k/\varepsilon^2)^{O(\sqrt{k})})$ . This is still exponential in the input size if  $k$  is large (say  $\sqrt{n}$  or  $n/100$ ), but if  $k$  and  $\varepsilon$  are fixed constants, the running time is linear in the number of input points.

### \*30.10 An FPTAS for Subset Sum

An approximation scheme whose running time, for any fixed  $\varepsilon$ , is polynomial in  $n$  is called a *polynomial-time approximation scheme* or *PTAS* (usually pronounced “pee taz”). If in addition the running time depends only polynomially on  $\varepsilon$ , the algorithm is called a *fully polynomial-time approximation scheme* or *FPTAS* (usually pronounced “eff pee taz”). For example, an approximation scheme with running time  $O(n^2/\varepsilon^2)$  is an FPTAS; an approximation scheme with running time  $O(n^{1/\varepsilon^6})$  is a PTAS but not an FPTAS; and our approximation scheme for  $k$ -center clustering is not a PTAS.

The last problem we'll consider is the *SUBSETSUM* problem: Given a set  $X$  containing  $n$  positive integers and a target integer  $t$ , determine whether  $X$  has a subset whose elements sum to  $t$ . The lecture notes on NP-completeness include a proof that *SUBSETSUM* is NP-hard. As stated, this problem doesn't allow any sort of approximation—the answer is either TRUE or FALSE.<sup>8</sup> So we will consider a related optimization problem instead: Given set  $X$  and integer  $t$ , find the subset of  $X$  whose sum is as large as possible but no larger than  $t$ .

<sup>7</sup>R. Z. Hwang, R. C. T. Lee, and R. C. Chan. The slab dividing approach to solve the Euclidean  $p$ -center problem. *Algorithmica* 9(1):1–22, 1993.

<sup>8</sup>Do, or do not. There is no ‘try’. (Are old one thousand when years you, alphabetical also in order talk will you.)

We have already seen a dynamic programming algorithm to solve the decision version SUBSETSUM in time  $O(nt)$ ; a similar algorithm solves the optimization version in the same time bound. Here is a different algorithm, whose running time does not depend on  $t$ :

```
SUBSETSUM( $X[1..n], t$ ):
   $S_0 \leftarrow \{0\}$ 
  for  $i \leftarrow 1$  to  $n$ 
     $S_i \leftarrow S_{i-1} \cup (S_{i-1} + X[i])$ 
    remove all elements of  $S_i$  bigger than  $t$ 
  return  $\max S_n$ 
```

Here  $S_{i-1} + X[i]$  denotes the set  $\{s + X[i] \mid s \in S_{i-1}\}$ . If we store each  $S_i$  in a sorted array, the  $i$ th iteration of the for-loop requires time  $O(|S_{i-1}|)$ . Each set  $S_i$  contains all possible subset sums for the first  $i$  elements of  $X$ ; thus,  $S_i$  has at most  $2^i$  elements. On the other hand, since every element of  $S_i$  is an integer between 0 and  $t$ , we also have  $|S_i| \leq t + 1$ . It follows that the total running time of this algorithm is  $\sum_{i=1}^n O(|S_{i-1}|) = O(\min\{2^n, nt\})$ .

Of course, this is only an estimate of worst-case behavior. If several subsets of  $X$  have the same sum, the sets  $S_i$  will have fewer elements, and the algorithm will be faster. The key idea for finding an approximate solution quickly is to ‘merge’ nearby elements of  $S_i$ —if two subset sums are nearly equal, ignore one of them. On the one hand, merging similar subset sums will introduce some error into the output, but hopefully not too much. On the other hand, by reducing the size of the set of sums we need to maintain, we will make the algorithm faster, hopefully significantly so.

Here is our approximation algorithm. We make only two changes to the exact algorithm: an initial sorting phase and an extra FILTERING step inside the main loop.

<b>FILTER(<math>Z[1..k], \delta</math>):</b> <pre>   SORT(Z)   <math>j \leftarrow 1</math>   <math>Y[j] \leftarrow Z[i]</math>   for <math>i \leftarrow 2</math> to <math>k</math>     if <math>Z[i] &gt; (1 + \delta) \cdot Y[j]</math>       <math>j \leftarrow j + 1</math>       <math>Y[j] \leftarrow Z[i]</math>   return <math>Y[1..j]</math></pre>	<b>APPROXSUBSETSUM(<math>X[1..n], k, \varepsilon</math>):</b> <pre>   SORT(X)   <math>R_0 \leftarrow \{0\}</math>   for <math>i \leftarrow 1</math> to <math>n</math>     <math>R_i \leftarrow R_{i-1} \cup (R_{i-1} + X[i])</math>     <math>R_i \leftarrow \text{FILTER}(R_i, \varepsilon/2n)</math>     remove all elements of <math>R_i</math> bigger than <math>t</math>   return <math>\max R_n</math></pre>
--	--

**Theorem 8.** APPROXSUBSETSUM returns a  $(1 + \varepsilon)$ -approximation of the optimal subset sum, given any  $\varepsilon$  such that  $0 < \varepsilon \leq 1$ .

**Proof:** The theorem follows from the following claim, which we prove by induction:

For any element  $s \in S_i$ , there is an element  $r \in R_i$  such that  $r \leq s \leq r \cdot (1 + \varepsilon n/2)^i$ .

The claim is trivial for  $i = 0$ . Let  $s$  be an arbitrary element of  $S_i$ , for some  $i > 0$ . There are two cases to consider: either  $s \in S_{i-1}$ , or  $s \in S_{i-1} + x_i$ .

- (1) Suppose  $s \in S_{i-1}$ . By the inductive hypothesis, there is an element  $r' \in R_{i-1}$  such that  $r' \leq s \leq r' \cdot (1 + \varepsilon n/2)^{i-1}$ . If  $r' \in R_i$ , the claim obviously holds. On the other hand, if  $r' \notin R_i$ , there must be an element  $r \in R_i$  such that  $r < r' \leq r(1 + \varepsilon n/2)$ , which implies that

$$r < r' \leq s \leq r' \cdot (1 + \varepsilon n/2)^{i-1} \leq r \cdot (1 + \varepsilon n/2)^i,$$

so the claim holds.

- (2) Suppose  $s \in S_{i-1} + x_i$ . By the inductive hypothesis, there is an element  $r' \in R_{i-1}$  such that  $r' \leq s - x_i \leq r' \cdot (1 + \varepsilon n/2)^{i-1}$ . If  $r' + x_i \in R_i$ , the claim obviously holds. On the other hand, if  $r' + x_i \notin R_i$ , there must be an element  $r \in R_i$  such that  $r < r' + x_i \leq r(1 + \varepsilon n/2)$ , which implies that

$$\begin{aligned} r < r' + x_i &\leq s \leq r' \cdot (1 + \varepsilon n/2)^{i-1} + x_i \\ &\leq (r - x_i) \cdot (1 + \varepsilon n/2)^i + x_i \\ &\leq r \cdot (1 + \varepsilon n/2)^i - x_i \cdot ((1 + \varepsilon n/2)^i - 1) \\ &\leq r \cdot (1 + \varepsilon n/2)^i. \end{aligned}$$

so the claim holds.

Now let  $s^* = \max S_n$  and  $r^* = \max R_n$ . Clearly  $r^* \leq s^*$ , since  $R_n \subseteq S_n$ . Our claim implies that there is some  $r \in R_n$  such that  $s^* \leq r \cdot (1 + \varepsilon/2n)^n$ . But  $r$  cannot be bigger than  $r^*$ , so  $s^* \leq r^* \cdot (1 + \varepsilon/2n)^n$ . The inequalities  $e^x \geq 1 + x$  for all  $x$ , and  $e^x \leq 2x + 1$  for all  $0 \leq x \leq 1$ , imply that  $(1 + \varepsilon/2n)^n \leq e^{\varepsilon/2} \leq 1 + \varepsilon$ .  $\square$

**Theorem 9.** APPROXSUBSETSUM runs in  $O((n^3 \log n)/\varepsilon)$  time.

**Proof:** Assuming we keep each set  $R_i$  in a sorted array, we can merge the two sorted arrays  $R_{i-1}$  and  $R_{i-1} + x_i$  in  $O(|R_{i-1}|)$  time. FILTERIN  $R_i$  and removing elements larger than  $t$  also requires only  $O(|R_{i-1}|)$  time. Thus, the overall running time of our algorithm is  $O(\sum_i |R_i|)$ ; to express this in terms of  $n$  and  $\varepsilon$ , we need to prove an upper bound on the size of each set  $R_i$ .

Let  $\delta = \varepsilon/2n$ . Because we consider the elements of  $X$  in increasing order, every element of  $R_i$  is between 0 and  $i \cdot x_i$ . In particular, every element of  $R_{i-1} + x_i$  is between  $x_i$  and  $i \cdot x_i$ . After FILTERING, at most one element  $r \in R_i$  lies in the range  $(1 + \delta)^k \leq r < (1 + \delta)^{k+1}$ , for any  $k$ . Thus, at most  $\lceil \log_{1+\delta} i \rceil$  elements of  $R_{i-1} + x_i$  survive the call to FILTER. It follows that

$$\begin{aligned} |R_i| &= |R_{i-1}| + \left\lceil \frac{\log i}{\log(1 + \delta)} \right\rceil \\ &\leq |R_{i-1}| + \left\lceil \frac{\log n}{\log(1 + \delta)} \right\rceil && [i \leq n] \\ &\leq |R_{i-1}| + \left\lceil \frac{2 \ln n}{\delta} \right\rceil && [e^x \leq 1 + 2x \text{ for all } 0 \leq x \leq 1] \\ &\leq |R_{i-1}| + \left\lceil \frac{n \ln n}{\varepsilon} \right\rceil && [\delta = \varepsilon/2n] \end{aligned}$$

Unrolling this recurrence into a summation gives us the upper bound  $|R_i| \leq i \cdot \lceil (n \ln n)/\varepsilon \rceil = O((n^2 \log n)/\varepsilon)$ .

We conclude that the overall running time of APPROXSUBSETSUM is  $O((n^3 \log n)/\varepsilon)$ , as claimed.  $\square$

## Exercises

1. (a) Prove that for any set of jobs, the makespan of the greedy assignment is at most  $(2 - 1/m)$  times the makespan of the optimal assignment, where  $m$  is the number of machines.  
 (b) Describe a set of jobs such that the makespan of the greedy assignment is exactly  $(2 - 1/m)$  times the makespan of the optimal assignment, where  $m$  is the number of machines.  
 (c) Describe an efficient algorithm to solve the minimum makespan scheduling problem *exactly* if every processing time  $T[i]$  is a power of two.
2. (a) Find the smallest graph (minimum number of edges) for which GREEDYVERTEXCOVER does not return the smallest vertex cover.  
 (b) For any integer  $n$ , describe an  $n$ -vertex graph for which GREEDYVERTEXCOVER returns a vertex cover of size  $OPT \cdot \Omega(\log n)$ .
3. (a) Find the smallest graph (minimum number of edges) for which DUMBVERTEXCOVER does not return the smallest vertex cover.  
 (b) Describe an infinite family of graphs for which DUMBVERTEXCOVER returns a vertex cover of size  $2 \cdot OPT$ .
4. Consider the following heuristic for constructing a vertex cover of a connected graph  $G$ : return the set of non-leaf nodes in any depth-first spanning tree of  $G$ .
  - (a) Prove that this heuristic returns a vertex cover of  $G$ .
  - (b) Prove that this heuristic returns a 2-approximation to the minimum vertex cover of  $G$ .
  - (c) Describe an infinite family of graphs for which this heuristic returns a vertex cover of size  $2 \cdot OPT$ .
5. Consider the following optimization version of the PARTITION problem. Given a set  $X$  of positive integers, our task is to partition  $X$  into disjoint subsets  $A$  and  $B$  such that  $\max\{\sum A, \sum B\}$  is as small as possible. This problem is clearly NP-hard. Determine the approximation ratio of the following polynomial-time approximation algorithm. Prove your answer is correct.

```

PARTITION( $X[1..n]$ ):
  Sort  $X$  in increasing order
   $a \leftarrow 0$ ;  $b \leftarrow 0$ 
  for  $i \leftarrow 1$  to  $n$ 
    if  $a < b$ 
       $a \leftarrow a + X[i]$ 
    else
       $b \leftarrow b + X[i]$ 
  return  $\max\{a, b\}$ 

```

6. The *chromatic number*  $\chi(G)$  of a graph  $G$  is the minimum number of colors required to color the vertices of the graph, so that every edge has endpoints with different colors. Computing the chromatic number exactly is NP-hard.

Prove that the following problem is also NP-hard: Given an  $n$ -vertex graph  $G$ , return any integer between  $\chi(G)$  and  $\chi(G) + 573$ . [Note: This does not contradict the possibility of a constant factor approximation algorithm.]

7. Let  $G = (V, E)$  be an undirected graph, each of whose vertices is colored either red, green, or blue. An edge in  $G$  is *boring* if its endpoints have the same color, and *interesting* if its endpoints have different colors. The *most interesting 3-coloring* is the 3-coloring with the maximum number of interesting edges, or equivalently, with the fewest boring edges. Computing the most interesting 3-coloring is NP-hard, because the standard 3-coloring problem is a special case.

- (a) Let  $zzz(G)$  denote the number of boring edges in the most interesting 3-coloring of a graph  $G$ . Prove that it is NP-hard to approximate  $zzz(G)$  within a factor of  $10^{100}$ .
- (b) Let  $wow(G)$  denote the number of interesting edges in the most interesting 3-coloring of  $G$ . Suppose we assign each vertex in  $G$  a *random* color from the set {red, green, blue}. Prove that the expected number of interesting edges is at least  $\frac{2}{3}wow(G)$ .

8. Consider the following algorithm for coloring a graph  $G$ .

```

TREECOLOR( $G$ ):
     $T \leftarrow$  any spanning tree of  $G$ 
    Color the tree  $T$  with two colors
     $c \leftarrow 2$ 
    for each edge  $(u, v) \in G \setminus T$ 
         $T \leftarrow T \cup \{(u, v)\}$ 
        if  $color(u) = color(v)$    «Try recoloring  $u$  with an existing color»
            for  $i \leftarrow 1$  to  $c$ 
                if no neighbor of  $u$  in  $T$  has color  $i$ 
                     $color(u) \leftarrow i$ 
        if  $color(u) = color(v)$    «Try recoloring  $v$  with an existing color»
            for  $i \leftarrow 1$  to  $c$ 
                if no neighbor of  $v$  in  $T$  has color  $i$ 
                     $color(v) \leftarrow i$ 
        if  $color(u) = color(v)$    «Give up and create a new color»
             $c \leftarrow c + 1$ 
             $color(u) \leftarrow c$ 

```

- (a) Prove that this algorithm colors any bipartite graph with just two colors.
  - (b) Let  $\Delta(G)$  denote the maximum degree of any vertex in  $G$ . Prove that this algorithm colors any graph  $G$  with at most  $\Delta(G)$  colors. This trivially implies that TREECOLOR is a  $\Delta(G)$ -approximation algorithm.
  - (c) Prove that TREECOLOR is *not* a constant-factor approximation algorithm.
9. The KNAPSACK problem can be defined as follows. We are given a finite set of elements  $X$  where each element  $x \in X$  has a non-negative *size* and a non-negative *value*, along with an integer *capacity*  $c$ . Our task is to determine the maximum total value among all subsets of  $X$  whose total size is at most  $c$ . This problem is NP-hard. Specifically, the optimization version of SUBSETSUM is a special case, where each element's value is equal to its size.

Determine the approximation ratio of the following polynomial-time approximation algorithm. Prove your answer is correct.

```
APPROXKNAPSACK( $X, c$ ):
    return max{GREEDYKNAPSACK( $X, c$ ), PICKBESTONE( $X, c$ )}
```

```
GREEDYKNAPSACK( $X, c$ ):
```

```
Sort  $X$  in decreasing order by the ratio  $value/size$ 
 $S \leftarrow 0; V \leftarrow 0$ 
for  $i \leftarrow 1$  to  $n$ 
    if  $S + size(x_i) > c$ 
        return  $V$ 
     $S \leftarrow S + size(x_i)$ 
     $V \leftarrow V + value(x_i)$ 
return  $V$ 
```

```
PICKBESTONE( $X, c$ ):
```

```
Sort  $X$  in increasing order by  $size$ 
 $V \leftarrow 0$ 
for  $i \leftarrow 1$  to  $n$ 
    if  $size(x_i) > c$ 
        return  $V$ 
    if  $value(x_i) > V$ 
         $V \leftarrow value(x_i)$ 
return  $V$ 
```

10. In the *bin packing* problem, we are given a set of  $n$  items, each with weight between 0 and 1, and we are asked to load the items into as few bins as possible, such that the total weight in each bin is at most 1. It's not hard to show that this problem is NP-Hard; this question asks you to analyze a few common approximation algorithms. In each case, the input is an array  $W[1..n]$  of weights, and the output is the number of bins used.

```
NEXTFIT( $W[1..n]$ ):
```

```
 $b \leftarrow 0$ 
 $Total[0] \leftarrow \infty$ 
for  $i \leftarrow 1$  to  $n$ 
    if  $Total[b] + W[i] > 1$ 
         $b \leftarrow b + 1$ 
         $Total[b] \leftarrow W[i]$ 
    else
         $Total[b] \leftarrow Total[b] + W[i]$ 
return  $b$ 
```

```
FIRSTFIT( $W[1..n]$ ):
```

```
 $b \leftarrow 0$ 
for  $i \leftarrow 1$  to  $n$ 
     $j \leftarrow 1; found \leftarrow \text{FALSE}$ 
    while  $j \leq b$  and  $found = \text{FALSE}$ 
        if  $Total[j] + W[i] \leq 1$ 
             $Total[j] \leftarrow Total[j] + W[i]$ 
             $found \leftarrow \text{TRUE}$ 
         $j \leftarrow j + 1$ 
    if  $found = \text{FALSE}$ 
         $b \leftarrow b + 1$ 
         $Total[b] = W[i]$ 
return  $b$ 
```

- (a) Prove that NEXTFIT uses at most twice the optimal number of bins.
- (b) Prove that FIRSTFIT uses at most twice the optimal number of bins.
- \*(c) Prove that if the weight array  $W$  is initially sorted in decreasing order, then FIRSTFIT uses at most  $(4 \cdot OPT + 1)/3$  bins, where  $OPT$  is the optimal number of bins. The following facts may be useful (but you need to prove them if your proof uses them):
  - In the packing computed by FIRSTFIT, every item with weight more than  $1/3$  is placed in one of the first  $OPT$  bins.
  - FIRSTFIT places at most  $OPT - 1$  items outside the first  $OPT$  bins.

11. Given a graph  $G$  with edge weights and an integer  $k$ , suppose we wish to partition the vertices of  $G$  into  $k$  subsets  $S_1, S_2, \dots, S_k$  so that the sum of the weights of the edges that cross the partition (that is, have endpoints in different subsets) is as large as possible.

- (a) Describe an efficient  $(1 - 1/k)$ -approximation algorithm for this problem.
- (b) Now suppose we wish to minimize the sum of the weights of edges that do *not* cross the partition. What approximation ratio does your algorithm from part (a) achieve for the new problem? Justify your answer.
12. The lecture notes describe a  $(3/2)$ -approximation algorithm for the metric traveling salesman problem. Here, we consider computing minimum-cost Hamiltonian *paths*. Our input consists of a graph  $G$  whose edges have weights that satisfy the triangle inequality. Depending upon the problem, we are also given zero, one, or two endpoints.
- (a) If our input includes zero endpoints, describe a  $(3/2)$ -approximation to the problem of computing a minimum cost Hamiltonian path.
  - (b) If our input includes one endpoint  $u$ , describe a  $(3/2)$ -approximation to the problem of computing a minimum cost Hamiltonian path that starts at  $u$ .
  - (c) If our input includes two endpoints  $u$  and  $v$ , describe a  $(5/3)$ -approximation to the problem of computing a minimum cost Hamiltonian path that starts at  $u$  and ends at  $v$ .
13. Suppose we are given a collection of  $n$  *jobs* to execute on a machine containing a row of  $m$  processors. When the  $i$ th job is executed, it occupies a *contiguous* set of  $\text{prox}[i]$  processors for  $\text{time}[i]$  seconds. A *schedule* for a set of jobs assigns each job an interval of processors and a starting time, so that no processor works on more than one job at any time. The *makespan* of a schedule is the time from the start to the finish of all jobs. Finally, the *parallel scheduling problem* asks us to compute the schedule with the smallest possible makespan.
- (a) Prove that the parallel scheduling problem is NP-hard.
  - (b) Give an algorithm that computes a 3-approximation of the minimum makespan of a set of jobs in  $O(m \log m)$  time. That is, if the minimum makespan is  $M$ , your algorithm should compute a schedule with make-span at most  $3M$ . You can assume that  $n$  is a power of 2.
14. Consider the greedy algorithm for metric TSP: start at an arbitrary vertex  $u$ , and at each step, travel to the closest unvisited vertex.
- (a) Show that the greedy algorithm for metric TSP is an  $O(\log n)$ -approximation, where  $n$  is the number of vertices. [*Hint: Argue that the  $k$ th least expensive edge in the tour output by the greedy algorithm has weight at most  $OPT/(n - k + 1)$ ; try  $k = 1$  and  $k = 2$  first.*]
  - \*(b) Show that the greedy algorithm for metric TSP is no better than an  $O(\log n)$ -approximation. That is, describe an infinite family of weighted graphs such that the greedy algorithm returns a cycle whose weight is  $\Omega(\log n)$  times the optimal TSP tour.

*Jeder Genießende meint, dem Baume habe es an der Frucht gelegen;  
aber ihm lag am Samen.  
[Everyone who enjoys thinks that the fundamental thing about trees is the fruit,  
but in fact it is the seed.]*

— Friedrich Wilhelm Nietzsche,  
*Vermischte Meinungen und Sprüche* [Mixed Opinions and Maxims] (1879)

*In view of all the deadly computer viruses that have been spreading lately,  
Weekend Update would like to remind you:  
When you link up to another computer,  
you're linking up to every computer that that computer has ever linked up to.*

— Dennis Miller, "Saturday Night Live", (c. 1985)

*Anything that, in happening, causes itself to happen again, happens again.*

— Douglas Adams (2005)

# Proof by Induction

Induction is a method for proving universally quantified propositions—statements about *all* elements of a (usually infinite) set. Induction is also the single most useful tool for reasoning about, developing, and analyzing algorithms. These notes give several examples of inductive proofs, along with a standard boilerplate and some motivation to justify (and help you remember) why induction works.

## 1 Prime Divisors: Proof by Smallest Counterexample

A **divisor** of a positive integer  $n$  is a positive integer  $p$  such that the ratio  $n/p$  is an integer. The integer 1 is a divisor of every positive integer (because  $n/1 = n$ ), and every integer is a divisor of itself (because  $n/n = 1$ ). A **proper divisor** of  $n$  is any divisor of  $n$  other than  $n$  itself. A positive integer is **prime** if it has *exactly* two divisors, which must be 1 and itself; equivalently, a number is prime if and only if 1 is its only proper divisor. A positive integer is **composite** if it has more than two divisors (or equivalently, more than one proper divisor). The integer 1 is neither prime nor composite, because it has exactly one divisor, namely itself.

Let's prove our first theorem:

**Theorem 1.** *Every integer greater than 1 has a prime divisor.*

The very first thing that you should notice, after reading just one word of the theorem, is that this theorem is *universally quantified*—it's a statement about *all* the elements of a set, namely, the set of positive integers larger than 1. If we were forced at gunpoint to write this sentence using fancy logic notation, the first character would be the universal quantifier  $\forall$ , pronounced ‘for all’. Fortunately, that won't be necessary.

There are only two ways to prove a universally quantified statement: directly or by contradiction. Let's say that again, louder: ***There are only two ways to prove a universally quantified statement: directly or by contradiction.*** Here are the standard templates for these two methods, applied to Theorem 1:

**Direct proof:** Let  $n$  be an arbitrary integer greater than 1.

... blah blah blah ...

Thus,  $n$  has at least one prime divisor.  $\square$

**Proof by contradiction:** For the sake of argument, assume there is an integer greater than 1 with no prime divisor.

Let  $n$  be an arbitrary integer greater than 1 with no prime divisor.

... blah blah blah ...

But that's just silly. Our assumption must be incorrect.  $\square$

The shaded boxes `... blah blah blah ...` indicate missing proof details (which you have to fill in). Most people find proofs by contradiction easier to discover than direct proofs, so let's try that first.

**Proof by contradiction:** For the sake of argument, assume there is an integer greater than 1 with no prime divisor.

Let  $n$  be an arbitrary integer greater than 1 with no prime divisor.

Since  $n$  is a divisor of  $n$ , and  $n$  has no prime divisors,  $n$  cannot be prime.

Thus,  $n$  must have at least one divisor  $d$  such that  $1 < d < n$ .

Let  $d$  be an arbitrary divisor of  $n$  such that  $1 < d < n$ .

Since  $n$  has no prime divisors,  $d$  cannot be prime.

Thus,  $d$  has at least one divisor  $d'$  such that  $1 < d' < d$ .

Let  $d'$  be an arbitrary divisor of  $d$  such that  $1 < d' < d$ .

Because  $d/d'$  is an integer,  $n/d' = (n/d) \cdot (d/d')$  is also an integer.

Thus,  $d'$  is also a divisor of  $n$ .

Since  $n$  has no prime divisors,  $d'$  cannot be prime.

Thus,  $d'$  has at least one divisor  $d''$  such that  $1 < d'' < d'$ .

Let  $d''$  be an arbitrary divisor of  $d'$  such that  $1 < d'' < d'$ .

Because  $d'/d''$  is an integer,  $n/d'' = (n/d') \cdot (d'/d'')$  is also an integer.

Thus,  $d''$  is also a divisor of  $n$ .

Since  $n$  has no prime divisors,  $d''$  cannot be prime.

... blah HELP! blah I'M STUCK IN AN INFINITE LOOP! blah ...

But that's just silly. Our assumption must be incorrect.  $\square$

We seem to be stuck in an infinite loop, looking at smaller and smaller divisors  $d > d' > d'' > \dots$ , none of which are prime. But this loop can't really be infinite. There are only  $n - 1$  positive integers smaller than  $n$ , so the proof *must* end after *at most*  $n - 1$  iterations. But how do we turn this observation into a formal proof? We need a single, self-contained proof for all integers  $n$ ; we're not allowed to write longer proofs for bigger integers. The trick is to jump directly to the *smallest* counterexample.

**Proof by smallest counterexample:** For the sake of argument, assume that there is an integer greater than 1 with no prime divisor.

Let  $n$  be **the smallest** integer greater than 1 with no prime divisor.

Since  $n$  is a divisor of  $n$ , and  $n$  has no prime divisors,  $n$  cannot be prime.

Thus,  $n$  has a divisor  $d$  such that  $1 < d < n$ .

Let  $d$  be a divisor of  $n$  such that  $1 < d < n$ .

**Because  $n$  is the smallest counterexample**,  $d$  has a prime divisor.

Let  $p$  be a prime divisor of  $d$ .

Because  $d/p$  is an integer,  $n/p = (n/d) \cdot (d/p)$  is also an integer.

Thus,  $p$  is also a divisor of  $n$ .

But this contradicts our assumption that  $n$  has no prime divisors!

So our assumption must be incorrect.  $\square$

Hooray, our first proof! We're done!

Um... well... no, we're definitely *not* done. That's a first draft up there, not a final polished proof. We don't write proofs just to convince ourselves; proofs are primarily a tool to convince other people. (In particular, 'other people' includes the people grading your homeworks and exams.) And while

proofs by contradiction are usually easier to *write*, direct proofs are almost always easier to *read*. So as a service to our audience (and our grade), let's transform our minimal-counterexample proof into a direct proof.

Let's first rewrite the indirect proof slightly, to make the structure more apparent. First, we break the assumption that  $n$  is the smallest counterexample into three simpler assumptions: (1)  $n$  is an integer greater than 1; (2)  $n$  has no prime divisors; and (3) there are no smaller counterexamples. Second, instead of dismissing the possibility that  $n$  is prime out of hand, we include an explicit case analysis.

**Proof by smallest counterexample:** Let  $n$  be an arbitrary integer greater than 1.

For the sake of argument, suppose  $n$  has no prime divisor.

**Assume that every integer  $k$  such that  $1 < k < n$  has a prime divisor.**

There are two cases to consider: Either  $n$  is prime, or  $n$  is composite.

- Suppose  $n$  is prime.

Then  $n$  is a prime divisor of  $n$ .

- Suppose  $n$  is composite.

Then  $n$  has a divisor  $d$  such that  $1 < d < n$ .

Let  $d$  be a divisor of  $n$  such that  $1 < d < n$ .

**Because no counterexample is smaller than  $n$ ,  $d$  has a prime divisor.**

Let  $p$  be a prime divisor of  $d$ .

Because  $d/p$  is an integer,  $n/p = (n/d) \cdot (d/p)$  is also an integer.

Thus,  $p$  is a prime divisor of  $n$ .

In each case, we conclude that  $n$  has a prime divisor.

But this contradicts our assumption that  $n$  has no prime divisors!

So our assumption must be incorrect. □

Now let's look carefully at the structure of this proof. First, we assumed that the statement we want to prove is false. Second, we proved that the statement we want to prove is true. Finally, we concluded from the contradiction that our assumption that the statement we want to prove is false is incorrect, so the statement we want to prove must be true.

But that's just silly. Why do we need the first and third steps? After all, the second step is a proof all by itself! Unfortunately, this redundant style of proof by contradiction is *extremely* common, even in professional papers. Fortunately, it's also very easy to avoid; just remove the first and third steps!

**Proof by induction:** Let  $n$  be an arbitrary integer greater than 1.

**Assume that every integer  $k$  such that  $1 < k < n$  has a prime divisor.**

There are two cases to consider: Either  $n$  is prime or  $n$  is composite.

- First, suppose  $n$  is prime.

Then  $n$  is a prime divisor of  $n$ .

- Now suppose  $n$  is composite.

Then  $n$  has a divisor  $d$  such that  $1 < d < n$ .

Let  $d$  be a divisor of  $n$  such that  $1 < d < n$ .

**Because no counterexample is smaller than  $n$ ,  $d$  has a prime divisor.**

Let  $p$  be a prime divisor of  $d$ .

Because  $d/p$  is an integer,  $n/p = (n/d) \cdot (d/p)$  is also an integer.

Thus,  $p$  is a prime divisor of  $n$ .

In both cases, we conclude that  $n$  has a prime divisor. □

This style of proof is called *induction*.<sup>1</sup> The assumption that there are no counterexamples smaller than  $n$  is called the *induction hypothesis*. The two cases of the proof have different names. The

<sup>1</sup>Many authors use the high-falutin' name *the principle of mathematical induction*, to distinguish it from *inductive reasoning*, the informal process by which we conclude that pigs can't whistle, horses can't fly, and NP-hard problems cannot be solved in polynomial time. We already know that every proof is mathematical (and arguably, all mathematics is proof), so as a description of a *proof* technique, the adjective 'mathematical' is simply redundant.

first case, which we argue directly, is called the *base case*. The second case, which actually uses the induction hypothesis, is called the *inductive case*. You may find it helpful to actually label the induction hypothesis, the base case(s), and the inductive case(s) in your proof.

The following point cannot be emphasized enough: The only difference between a proof by induction and a proof by smallest counterexample is the way we write down the argument. The essential structure of the proofs are exactly the same. The core of our original indirect argument is a proof of the following implication for all  $n$ :

$$n \text{ has no prime divisor} \implies \text{some number smaller than } n \text{ has no prime divisor.}$$

The core of our direct proof is the following logically equivalent implication:

$$\text{every number smaller than } n \text{ has a prime divisor} \implies n \text{ has a prime divisor}$$

The left side of this implication is just the induction hypothesis.

The proofs we've been playing with have been very careful and explicit; until you're comfortable writing your own proofs, you should be equally careful. A more mature proof-writer might express the same proof more succinctly as follows:

**Proof by induction:** Let  $n$  be an arbitrary integer greater than 1. Assume that every integer  $k$  such that  $1 < k < n$  has a prime divisor. If  $n$  is prime, then  $n$  is a prime divisor of  $n$ . On the other hand, if  $n$  is composite, then  $n$  has a proper divisor  $d$ . The induction hypothesis implies that  $d$  has a prime divisor  $p$ . The integer  $p$  is also a divisor of  $n$ .  $\square$

A proof in this more succinct form is still worth full credit, provided the induction hypothesis is written explicitly and the case analysis is obviously exhaustive.

A professional mathematician would write the proof even more tersely:

**Proof:** Induction.  $\square$

And you can write that tersely, too, when you're a professional mathematician.

## 2 The Axiom of Induction

Why does this work? Well, let's step back to the original proof by smallest counterexample. How do we know that a smallest counterexample exists? This seems rather obvious, but in fact, it's impossible to prove without using the following fact:

*Every non-empty set of positive integers has a smallest element.*

Every set  $X$  of positive integers is the set of counterexamples to some proposition  $P(n)$  (specifically, the proposition  $n \notin X$ ). Thus, the Axiom of Induction can be rewritten as follows:

*If the proposition  $P(n)$  is false for some positive integer  $n$ ,  
then  
the proposition  $(P(1) \wedge P(2) \wedge \dots \wedge P(n-1) \wedge \neg P(n))$  is true for some positive integer  $n$ .*

Equivalently, in English:

*If some statement about positive integers has a counterexample,  
then  
that statement has a smallest counterexample.*

We can write this implication in contrapositive form as follows:

*If the proposition  $(P(1) \wedge P(2) \wedge \cdots \wedge P(n-1) \wedge \neg P(n))$  is false for every positive integer  $n$ ,  
then  
the proposition  $P(n)$  is true for every positive integer  $n$ .*

Finally, let's rewrite the first half of this statement in a logically equivalent form, by replacing  $\neg(p \wedge \neg q)$  with  $p \rightarrow q$ .

*If the implication  $(P(1) \wedge P(2) \wedge \cdots \wedge P(n-1)) \rightarrow P(n)$  is true for every positive integer  $n$ ,  
then  
the proposition  $P(n)$  is true for every positive integer  $n$ .*

This formulation is usually called the **Axiom of Induction**. In a proof by induction that  $P(n)$  holds for all  $n$ , the conjunction  $(P(1) \wedge P(2) \wedge \cdots \wedge P(n-1))$  is the inductive hypothesis.

A **proof by induction** for the proposition “ $P(n)$  for every positive integer  $n$ ” is nothing but a direct proof of the more complex proposition “ $(P(1) \wedge P(2) \wedge \cdots \wedge P(n-1)) \rightarrow P(n)$  for every positive integer  $n$ ”. Because it's a direct proof, it *must* start by considering an arbitrary positive integer, which we might as well call  $n$ . Then, to prove the implication, we explicitly assume the hypothesis  $(P(1) \wedge P(2) \wedge \cdots \wedge P(n-1))$  and then prove the conclusion  $P(n)$  for that particular value of  $n$ . The proof almost always breaks down into two or more cases, each of which may or may not actually use the inductive hypothesis.

Here is the boilerplate for every induction proof. Read it. Learn it. Use it.

**Theorem:**  $P(n)$  for every positive integer  $n$ .

**Proof by induction:** Let  $n$  be an arbitrary positive integer.

Assume inductively that  $P(k)$  is true for every positive integer  $k < n$ .

There are several cases to consider:

- Suppose  $n$  is ... blah blah blah ...  
Then  $P(n)$  is true.
- Suppose  $n$  is ... blah blah blah ...

The inductive hypothesis implies that ... blah blah blah ...

Thus,  $P(n)$  is true.

In each case, we conclude that  $P(n)$  is true. □

Some textbooks distinguish between several different types of induction: ‘regular’ induction versus ‘strong’ induction versus ‘complete’ induction versus ‘structural’ induction versus ‘transfinite’ induction versus ‘Noetherian’ induction. Distinguishing between these different types of induction is pointless hairsplitting; I won’t even define them. Every ‘different type’ of induction proof is provably equivalent to a proof by smallest counterexample. (Later we will consider inductive proofs of statements about partially ordered sets other than the positive integers, for which ‘smallest’ has a different meaning, but this difference will prove to be inconsequential.)

### 3 Stamps and Recursion

Let’s move on to a completely different example.

**Theorem 2.** Given an unlimited supply of 5-cent stamps and 7-cent stamps, we can make any amount of postage larger than 23 cents.

We could prove this by contradiction, using a smallest-counterexample argument, but let's aim for a direct proof by induction this time. We start by writing down the induction boilerplate, using the standard induction hypothesis: *There is no counterexample smaller than  $n$ .*

**Proof by induction:** Let  $n$  be an arbitrary integer greater than 23.

Assume that for any integer  $k$  such that  $23 < k < n$ , we can make  $k$  cents in postage.

... blah blah blah ...

Thus, we can make  $n$  cents in postage.  $\square$

How do we fill in the details? One approach is to think about what you would actually do if you really had to make  $n$  cents in postage. For example, you might start with a 5-cent stamp, and then try to make  $n - 5$  cents in postage. The inductive hypothesis says you can make *any* amount of postage bigger than 23 cents and less than  $n$  cents. So if  $n - 5 > 23$ , then you *already know* that you can make  $n - 5$  cents in postage! (You don't know *how* to make  $n - 5$  cents in postage, but so what?)

Let's write this observation into our proof as two separate cases: either  $n > 28$  (where our approach works) or  $n \leq 28$  (where we don't know what to do yet).

**Proof by induction:** Let  $n$  be an arbitrary integer greater than 23.

Assume that for any integer  $k$  such that  $23 < k < n$ , we can make  $k$  cents in postage.

There are two cases to consider: Either  $n > 28$  or  $n \leq 28$ .

- Suppose  $n > 28$ .

Then  $23 < n - 5 < n$ .

Thus, **the induction hypothesis** implies that we can make  $n - 5$  cents in postage.

Adding one more 5-cent stamp gives us  $n$  cents in postage.

- Now suppose  $n \leq 28$ .

... blah blah blah ...

In both cases, we can make  $n$  cents in postage.  $\square$

What do we do in the second case? Fortunately, this case considers only five integers: 24, 25, 26, 27, and 28. There might be a clever way to solve all five cases at once, but why bother? They're small enough that we can find a solution by brute force in less than a minute. To make the proof more readable, I'll unfold the nested cases and list them in increasing order.

**Proof by induction:** Let  $n$  be an arbitrary integer greater than 23.

Assume that for any integer  $k$  such that  $23 < k < n$ , we can make  $k$  cents in postage.

There are six cases to consider:  $n = 24$ ,  $n = 25$ ,  $n = 26$ ,  $n = 27$ ,  $n = 28$ , and  $n > 28$ .

- $24 = 7 + 7 + 5 + 5$
- $25 = 5 + 5 + 5 + 5 + 5$
- $26 = 7 + 7 + 7 + 5$
- $27 = 7 + 5 + 5 + 5 + 5$
- $28 = 7 + 7 + 7 + 7$
- Suppose  $n > 28$ .

Then  $23 < n - 5 < n$ .

Thus, **the induction hypothesis** implies that we can make  $n - 5$  cents in postage.

Adding one more 5-cent stamp gives us  $n$  cents in postage.

In all cases, we can make  $n$  cents in postage.  $\square$

Voilà! An induction proof! More importantly, we now have a recipe for *discovering* induction proofs.

1. **Write down the boilerplate.** Write down the universal invocation ('Let  $n$  be an arbitrary...'), the induction hypothesis, and the conclusion, with enough blank space for the remaining details. Don't be clever. Don't even think. Just write. ***This is the easy part.*** To emphasize the common structure, the boilerplate will be indicated in green for the rest of this handout.
2. **Think big.** Don't think how to solve the problem all the way down to the ground; you'll only make yourself dizzy. Don't think about piddly little numbers like 1 or 5 or  $10^{100}$ . Instead, think about how to reduce the proof about some *absoluckingly ginormous* value of  $n$  to a proof about some other number(s) smaller than  $n$ . ***This is the hard part.***
3. **Look for holes.** Look for cases where your inductive argument breaks down. Solve those cases directly. Don't be clever here; be stupid but thorough.
4. **Rewrite everything.** Your first proof is a rough draft. Rewrite the proof so that your argument is easier for your (unknown?) reader to follow.

The cases in an inductive proof always fall into two categories. Any case that uses the inductive hypothesis is called an *inductive case*. Any case that does not use the inductive hypothesis is called a *base case*. Typically, but *not* always, base cases consider a few small values of  $n$ , and the inductive cases consider everything else. Induction proofs are usually clearer if we present the base cases first, but I find it much easier to *discover* the inductive cases first. In other words, I recommend writing induction proofs backwards.

Well-written induction proofs very closely resemble well-written recursive programs. We computer scientists use induction primarily to reason about recursion, so maintaining this resemblance is extremely useful—we only have to keep one mental pattern, called ‘induction’ when we’re writing proofs and ‘recursion’ when we’re writing code. Consider the following C and Scheme programs for making  $n$  cents in postage:

```
void postage(int n)
{
    assert(n>23);
    switch ($n$)
    {
        case 24: printf("7+7+5+5"); break;
        case 25: printf("5+5+5+5+5"); break;
        case 26: printf("7+7+7+5"); break;
        case 27: printf("7+5+5+5+5"); break;
        case 28: printf("7+7+7+7"); break;
        default:
            postage(n-5);
            printf("+5");
    }
}
```

```
(define (postage n)
  (cond ((= n 24) (5 5 7 7))
        ((= n 25) (5 5 5 5 5))
        ((= n 26) (5 7 7 7))
        ((= n 27) (5 5 5 5 7))
        ((= n 28) (7 7 7 7))
        ((> n 28) (cons 5 (postage (- n 5))))))
```

The C program begins by declaring the input parameter (“Let  $n$  be an arbitrary integer...”) and asserting its range (“... greater than 23.”). (Scheme programs don’t have type declarations.) In both languages, the code branches into six cases: five that are solved directly, plus one that is handled by invoking the inductive hypothesis recursively.

## 4 More on Prime Divisors

Before we move on to different examples, let’s prove another fact about prime numbers:

**Theorem 3.** *Every positive integer is a product of prime numbers.*

First, let’s write down the boilerplate. Hey! I saw that! You were *thinking*, weren’t you? Stop that this instant! Don’t make me turn the car around. **First** we write down the boilerplate.

**Proof by induction:** Let  $n$  be an arbitrary positive integer.

Assume that any positive integer  $k < n$  is a product of prime numbers.

There are **some** cases to consider:

... blah blah blah ...

Thus,  $n$  is a product of prime numbers.  $\square$

Now let’s think about how you would actually factor a positive integer  $n$  into primes. There are a couple of different options here. One possibility is to find a prime divisor  $p$  of  $n$ , as guaranteed by Theorem 1, and recursively factor the integer  $n/p$ . This argument works as long as  $n \geq 2$ , but what about  $n = 1$ ? The answer is simple: 1 is the product of the *empty* set of primes. What else could it be?

**Proof by induction:** Let  $n$  be an arbitrary positive integer.

Assume that any positive integer  $k < n$  is a product of prime numbers.

There are two **cases to consider**: either  $n = 1$  or  $n \geq 2$ .

- If  $n = 1$ , then  $n$  is the product of the elements of the empty set, each of which is prime, green, sparkly, vanilla, and hemophagic.
- Suppose  $n > 1$ . Let  $p$  be a prime divisor of  $n$ , as guaranteed by Theorem 2. The inductive hypothesis implies that the positive integer  $n/p$  is a product of primes, and clearly  $n = (n/p) \cdot p$ .

In both cases,  $n$  is a product of prime numbers.  $\square$

But an even simpler method is to factor  $n$  into any two proper divisors, and recursively handle them both. This method works as long as  $n$  is composite, since otherwise there is no way to factor  $n$  into smaller integers. Thus, we need to consider prime numbers separately, as well as the special case 1.

**Proof by induction:** Let  $n$  be an arbitrary positive integer.

Assume that any positive integer  $k < n$  is a product of prime numbers.

There are three **cases to consider**: either  $n = 1$ ,  $n$  is prime, or  $n$  is composite.

- If  $n = 1$ , then  $n$  is the product of the elements of the empty set, each of which is prime, red, broody, chocolate, and lycanthropic.
- If  $n$  is prime, then  $n$  is the product of one prime number, namely  $n$ .
- Suppose  $n$  is composite. Let  $d$  be any proper divisor of  $n$  (guaranteed by the definition of ‘composite’), and let  $m = n/d$ . Since both  $d$  and  $m$  are positive integers smaller than  $n$ , the inductive hypothesis implies that  $d$  and  $m$  are both products of prime numbers. We clearly have  $n = d \cdot m$ .

In both cases,  $n$  is a product of prime numbers.  $\square$

## 5 Summations

Here's an easy one.

**Theorem 4.**  $\sum_{i=0}^n 3^i = \frac{3^{n+1} - 1}{2}$  for every non-negative integer  $n$ .

First let's write down the induction boilerplate, which empty space for the details we'll fill in later.

**Proof by induction:** Let  $n$  be an arbitrary non-negative integer.

Assume inductively that  $\sum_{i=0}^k 3^i = \frac{3^{k+1} - 1}{2}$  for every non-negative integer  $k < n$ .

There are *some number of* cases to consider:

... blah blah blah ...

We conclude that  $\sum_{i=0}^n 3^i = \frac{3^{n+1} - 1}{2}$ . □

Now imagine you are part of an infinitely long assembly line of mathematical provers, each assigned to a particular non-negative integer. Your task is to prove this theorem for the integer 8675310. The regulations of the Mathematical Provers Union require you not to think about any other integer but your own. The assembly line starts with the Senior Master Prover, who proves the theorem for the case  $n = 0$ . Next is the Assistant Senior Master Prover, who proves the theorem for  $n = 1$ . After him is the Assistant Assistant Senior Master Prover, who proves the theorem for  $n = 2$ . Then the Assistant Assistant Assistant Senior Master Prover proves the theorem for  $n = 3$ . As the work proceeds, you start to get more and more bored. You attempt strike up a conversation with Jenny, the prover to your left, but she ignores you, preferring to focus on the proof. Eventually, you fall into a deep, dreamless sleep. An undetermined time later, Jenny wakes you up by shouting, "Hey, doofus! It's your turn!" As you look around, bleary-eyed, you realize that Jenny and everyone to your left has finished their proofs, and that everyone is waiting for you to finish yours. What do you do?

What you do, after wiping the drool off your chin, is stop and think for a moment about what you're trying to prove. What does that  $\sum$  notation actually mean? Intuitively, we can expand the notation as follows:

$$\sum_{i=0}^{8675310} 3^i = 3^0 + 3^1 + \cdots + 3^{8675309} + 3^{8675310}.$$

Notice that this expression also contains the summation that Jenny just finished proving something about:

$$\sum_{i=0}^{8675309} 3^i = 3^0 + 3^1 + \cdots + 3^{8675308} + 3^{8675309}.$$

Putting these two expressions together gives us the following identity:

$$\sum_{i=0}^{8675310} 3^i = \sum_{i=0}^{8675309} 3^i + 3^{8675310}$$

In fact, this recursive identity is the *definition* of  $\sum$ . Jenny just proved that the summation on the right is equal to  $(3^{8675310} - 1)/2$ , so we can plug that into the right side of our equation:

$$\sum_{i=0}^{8675310} 3^i = \sum_{i=0}^{8675309} 3^i + 3^{8675310} = \frac{3^{8675310} - 1}{2} + 3^{8675310}.$$

And it's all downhill from here. After a little bit of algebra, you simplify the right side of this equation to  $(3^{8675311} - 1)/2$ , wake up the prover to your right, and start planning your well-earned vacation.

Let's insert this argument into our boilerplate, only using a generic 'big' integer  $n$  instead of the specific integer 8675310:

**Proof by induction:** Let  $n$  be an arbitrary non-negative integer.

Assume inductively that  $\sum_{i=0}^k 3^i = \frac{3^{k+1} - 1}{2}$  for every non-negative integer  $k < n$ .

There are two cases to consider: Either  $n$  is big or  $n$  is small.

- If  $n$  is big, then

$$\begin{aligned}\sum_{i=0}^n 3^i &= \sum_{i=0}^{n-1} 3^i + 3^n && [\text{definition of } \sum] \\ &= \frac{3^n - 1}{2} + 3^n && [\text{induction hypothesis, with } k = n - 1] \\ &= \frac{3^{n+1} - 1}{2} && [\text{algebra}]\end{aligned}$$

- On the other hand, if  $n$  is small, then ... blah blah blah ...

In both cases, we conclude that  $\sum_{i=0}^n 3^i = \frac{3^{n+1} - 1}{2}$ . □

Now, how big is 'big', and what do we do when  $n$  is 'small'? To answer the first question, let's look at where our existing inductive argument breaks down. In order to apply the induction hypothesis when  $k = n - 1$ , the integer  $n - 1$  must be non-negative; equivalently,  $n$  must be positive. But that's the only assumption we need: **The only case we missed is  $n = 0$** . Fortunately, this case is easy to handle directly.

**Proof by induction:** Let  $n$  be an arbitrary non-negative integer.

Assume inductively that  $\sum_{i=0}^k 3^i = \frac{3^{k+1} - 1}{2}$  for every non-negative integer  $k < n$ .

There are two cases to consider: Either  $n = 0$  or  $n \geq 1$ .

- If  $n = 0$ , then  $\sum_{i=0}^n 3^i = 3^0 = 1$ , and  $\frac{3^{n+1} - 1}{2} = \frac{3^1 - 1}{2} = 1$ .
- On the other hand, if  $n \geq 1$ , then

$$\begin{aligned}\sum_{i=0}^n 3^i &= \sum_{i=0}^{n-1} 3^i + 3^n && [\text{definition of } \sum] \\ &= \frac{3^n - 1}{2} + 3^n && [\text{induction hypothesis, with } k = n - 1] \\ &= \frac{3^{n+1} - 1}{2} && [\text{algebra}]\end{aligned}$$

In both cases, we conclude that  $\sum_{i=0}^n 3^i = \frac{3^{n+1} - 1}{2}$ . □

Here is the same proof, written more tersely; the non-standard symbol  $\stackrel{IH}{\equiv}$  indicates the use of the induction hypothesis.

**Proof by induction:** Let  $n$  be an arbitrary non-negative integer, and assume inductively that  $\sum_{i=0}^k 3^i = (3^{k+1} - 1)/2$  for every non-negative integer  $k < n$ . The base case  $n = 0$  is trivial, and for any  $n \geq 1$ , we have

$$\sum_{i=0}^n 3^i = \sum_{i=0}^{n-1} 3^i + 3^n \stackrel{IH}{\equiv} \frac{3^n - 1}{2} + 3^n = \frac{3^{n+1} - 1}{2}.$$

□

This is not the only way to prove this theorem by induction; here is another:

**Proof by induction:** Let  $n$  be an arbitrary non-negative integer, and assume inductively that  $\sum_{i=0}^k 3^i = (3^{k+1} - 1)/2$  for every non-negative integer  $k < n$ . The base case  $n = 0$  is trivial, and for any  $n \geq 1$ , we have

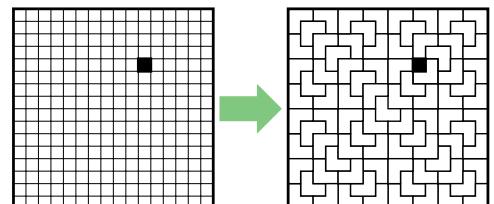
$$\sum_{i=0}^n 3^i = 3^0 + \sum_{i=1}^n 3^i = 3^0 + 3 \cdot \sum_{i=0}^{n-1} 3^i \stackrel{IH}{=} 3^0 + 3 \cdot \frac{3^n - 1}{2} = \frac{3^{n+1} - 1}{2}.$$

□

In the remainder of these notes, I'll give several more examples of induction proofs. In some cases, I give multiple proofs for the same theorem. Unlike the earlier examples, I will not describe the thought process that lead to the proof; in each case, I followed the basic outline on page 7.

## 6 Tiling with Triominos

The next theorem is about *tiling* a square checkerboard with *triominos*. A triomino is a shape composed of three squares meeting in an L-shape. Our goal is to cover as much of a  $2^n \times 2^n$  grid with triominos as possible, without any two triominos overlapping, and with all triominos inside the square. We can't cover every square in the grid—the number of squares is  $4^n$ , which is not a multiple of 3—but we can cover all but one square. In fact, as the next theorem shows, we can choose *any* square to be the one we don't want to cover.



Almost tiling a  $16 \times 16$  checkerboard with triominos.

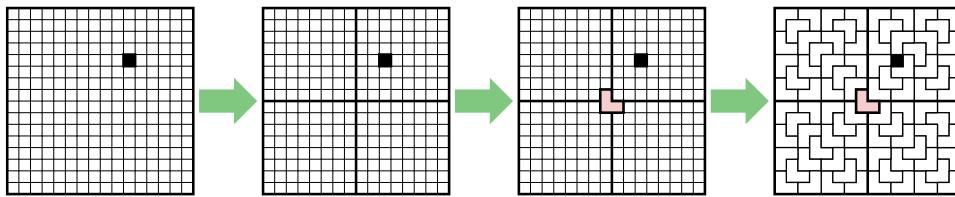
**Theorem 5.** For any non-negative integer  $n$ , the  $2^n \times 2^n$  checkerboard with any square removed can be tiled using L-shaped triominos.

Here are two different inductive proofs for this theorem, one ‘top down’, the other ‘bottom up’.

**Proof by top-down induction:** Let  $n$  be an arbitrary non-negative integer. Assume that for any non-negative integer  $k < n$ , the  $2^k \times 2^k$  grid with any square removed can be tiled using triominos. There are two cases to consider: Either  $n = 0$  or  $n \geq 1$ .

- The  $2^0 \times 2^0$  grid has a single square, so removing one square leaves nothing, which we can tile with zero triominos.
- Suppose  $n \geq 1$ . In this case, the  $2^n \times 2^n$  grid can be divided into four smaller  $2^{n-1} \times 2^{n-1}$  grids. Without loss of generality, suppose the deleted square is in the upper right quarter. With a single L-shaped triomino at the center of the board, we can cover one square in each of the other three quadrants. The induction hypothesis implies that we can tile each of the quadrants, minus one square.

In both cases, we conclude that the  $2^n \times 2^n$  grid with any square removed can be tiled with triominos.  $\square$

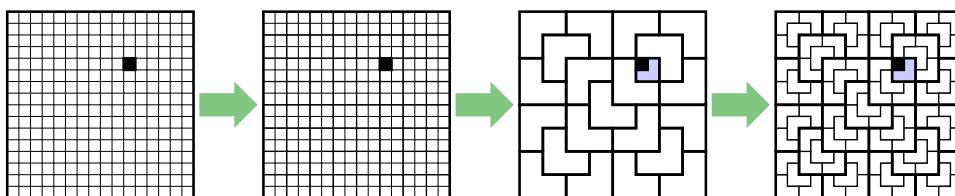


Top-down inductive proof of Theorem 4.

**Proof by bottom-up induction:** Let  $n$  be an arbitrary non-negative integer. Assume that for any non-negative integer  $k < n$ , the  $2^k \times 2^k$  grid with any square removed can be tiled using triominos. There are two cases to consider: Either  $n = 0$  or  $n \geq 1$ .

- The  $2^0 \times 2^0$  grid has a single square, so removing one square leaves nothing, which we can tile with zero triominos.
- Suppose  $n \geq 1$ . Then by clustering the squares into  $2 \times 2$  blocks, we can transform any  $2^n \times 2^n$  grid into a  $2^{n-1} \times 2^{n-1}$  grid. Suppose square  $(i, j)$  has been removed from the  $2^n \times 2^n$  grid. The induction hypothesis implies that the  $2^{n-1} \times 2^{n-1}$  grid with block  $([i/2], [j/2])$  removed can be tiled with double-size triominos. Each double-size triomino can be tiled with four smaller triominos, and block  $([i/2], [j/2])$  with square  $(i, j)$  removed is another triomino.

In both cases, we conclude that the  $2^n \times 2^n$  grid with any square removed can be tiled with triominos.  $\square$



Second proof of Theorem 4.

## 7 Binary Numbers Exist

**Theorem 6.** Every non-negative integer can be written as the sum of distinct powers of 2.

Intuitively, this theorem states that every number can be represented in binary. (That's not a proof, by the way; it's just a restatement of the theorem.) I'll present *four* distinct inductive proofs for this theorem. The first two are standard, by-the-book induction proofs.

**Proof by top-down induction:** Let  $n$  be an arbitrary non-negative integer. Assume that any non-negative integer less than  $n$  can be written as the sum of distinct powers of 2.

There are two cases to consider: Either  $n = 0$  or  $n \geq 1$ .

- The base case  $n = 0$  is trivial—the elements of the empty set are distinct and sum to zero.
- Suppose  $n \geq 1$ . Let  $k$  be the largest integer such that  $2^k \leq n$ , and let  $m = n - 2^k$ . Observe that  $m < 2^{k+1} - 2^k = 2^k$ . Because  $0 \leq m < n$ , the inductive hypothesis implies that  $m$  can be written as the sum of distinct powers of 2. Moreover, in the summation for  $m$ , each power of 2 is at most  $m$ , and therefore less than  $2^k$ . Thus,  $m + 2^k$  is the sum of distinct powers of 2.

In either case, we conclude that  $n$  can be written as the sum of distinct powers of 2.  $\square$

**Proof by bottom-up induction:** Let  $n$  be an arbitrary non-negative integer. Assume that any non-negative integer less than  $n$  can be written as the sum of distinct powers of 2. There are two cases to consider: Either  $n = 0$  or  $n \geq 1$ .

- The base case  $n = 0$  is trivial—the elements of the empty set are distinct and sum to zero.
- Suppose  $n \geq 1$ , and let  $m = \lfloor n/2 \rfloor$ . Because  $0 \leq m < n$ , the inductive hypothesis implies that  $m$  can be written as the sum of distinct powers of 2. Thus,  $2m$  can also be written as the sum of distinct powers of 2, each of which is greater than  $2^0$ . If  $n$  is even, then  $n = 2m$  and we are done; otherwise,  $n = 2m + 2^0$  is the sum of distinct powers of 2.

In either case, we conclude that  $n$  can be written as the sum of distinct powers of 2.  $\square$

The third proof deviates slightly from the induction boilerplate. At the top level, this proof doesn't actually use induction at all! However, a key step requires its own (straightforward) inductive proof.

**Proof by algorithm:** Let  $n$  be an arbitrary non-negative integer. Let  $S$  be a multiset containing  $n$  copies of  $2^0$ . Modify  $S$  by running the following algorithm:

```
while S has more than one copy of any element  $2^i$ 
    Remove two copies of  $2^i$  from S
    Insert one copy of  $2^{i+1}$  into S
```

Each iteration of this algorithm reduces the cardinality of  $S$  by 1, so the algorithm must eventually halt. When the algorithm halts, the elements of  $S$  are distinct. We claim that just after each iteration of the while loop, the elements of  $S$  sum to  $n$ .

**Proof by induction:** Consider an arbitrary iteration of the loop. Assume inductively that just after each previous iteration, the elements of  $S$  sum to  $n$ . Before any iterations of the loop, the elements of  $S$  sum to  $n$  by definition. The induction hypothesis implies that just before the current iteration begins, the elements of  $S$  sum to  $n$ . The loop replaces two copies of some number  $2^i$  with their sum  $2^{i+1}$ , leaving the total sum of  $S$  unchanged. Thus, when the iteration ends, the elements of  $S$  sum to  $n$ .  $\square$

Thus, when the algorithm halts, the elements of  $S$  are distinct powers of 2 that sum to  $n$ . We conclude that  $n$  can be written as the sum of distinct powers of 2.  $\square$

The fourth proof uses so-called ‘weak’ induction, where the inductive hypothesis can only be applied at  $n - 1$ . Not surprisingly, tying all but one hand behind our backs makes the resulting proof longer, more complicated, and harder to read. It doesn’t help that the algorithm used in the proof is overly specific. Nevertheless, this is the first approach that occurs to most students who have not truly accepted the Recursion Fairy into their hearts.

**Proof by baby-step induction:** Let  $n$  be an arbitrary non-negative integer. Assume that any non-negative integer less than  $n$  can be written as the sum of distinct powers of 2. There are two cases to consider: Either  $n = 0$  or  $n \geq 1$ .

- The base case  $n = 0$  is trivial—the elements of the empty set are distinct and sum to zero.
- Suppose  $n \geq 1$ . The inductive hypothesis implies that  $n - 1$  can be written as the sum of distinct powers of 2. Thus,  $n$  can be written as the sum of powers of 2, which are distinct except possibly for two copies of  $2^0$ . Let  $S$  be this multiset of powers of 2.

Now consider the following algorithm:

```
i ← 0
while S has more than one copy of 2i
    Remove two copies of 2i from S
    Insert one copy of 2i+1 into S
    i ← i + 1
```

Each iteration of this algorithm reduces the cardinality of  $S$  by 1, so the algorithm must eventually halt. We claim that for every non-negative integer  $i$ , the following invariants are satisfied after the  $i$ th iteration of the while loop (or before the algorithm starts if  $i = 0$ ):

- The elements of  $S$  sum to  $n$ .

**Proof by induction:** Let  $i$  be an arbitrary non-negative integer. Assume that for any non-negative integer  $j \leq i$ , after the  $j$ th iteration of the while loop, the elements of  $S$  sum to  $n$ . If  $i = 0$ , the elements of  $S$  sum to  $n$  by definition of  $S$ . Otherwise, the induction hypothesis implies that just before the  $i$ th iteration, the elements of  $S$  sum to  $n$ ; the  $i$ th iteration replaces two copies of  $2^i$  with  $2^{i+1}$ , leaving the sum unchanged.  $\square$

- The elements in  $S$  are distinct, except possibly for two copies of  $2^i$ .

**Proof by induction:** Let  $i$  be an arbitrary non-negative integer. Assume that for any non-negative integer  $j \leq i$ , after the  $j$ th iteration of the while loop, the elements of  $S$  are distinct except possibly for two copies of  $2^j$ . If  $i = 0$ , the invariant holds by definition of  $S$ . So suppose  $i > 0$ . The induction hypothesis implies that just before the  $i$ th iteration, the elements of  $S$  are distinct except possibly for two copies of  $2^i$ . If there are two copies of  $2^i$ , the algorithm replaces them both with  $2^{i+1}$ , and the invariant is established; otherwise, the algorithm halts, and the invariant is again established.  $\square$

The second invariant implies that when the algorithm halts, the elements of  $S$  are distinct.

In either case, we conclude that  $n$  can be written as the sum of distinct powers of 2.  $\square$

Repeat after me: “Doctor! Doctor! It hurts when I do this!”

## 8 Irrational Numbers Exist

**Theorem 7.**  $\sqrt{2}$  is irrational.

**Proof:** I will prove that  $p^2 \neq 2q^2$  (and thus  $p/q \neq \sqrt{2}$ ) for all positive integers  $p$  and  $q$ .

Let  $p$  and  $q$  be arbitrary positive integers. Assume that for any positive integers  $i < p$  and  $j < q$ , we have  $i^2 \neq 2j^2$ . Let  $i = \lfloor p/2 \rfloor$  and  $j = \lfloor q/2 \rfloor$ . There are three cases to consider:

- Suppose  $p$  is odd. Then  $p^2 = (2i + 1)^2 = 4i^2 + 4i + 1$  is odd, but  $2q^2$  is even.
- Suppose  $p$  is even and  $q$  is odd. Then  $p^2 = 4i^2$  is divisible by 4, but  $2q^2 = 2(2j+1)^2 = 4(2j^2 + 2j) + 2$  is not divisible by 4.
- Finally, suppose  $p$  and  $q$  are both even. The induction hypothesis implies that  $i^2 \neq 2j^2$ . Thus,  $p^2 = 4i^2 \neq 8j^2 = 2q^2$ .

In every case, we conclude that  $p^2 \neq 2q^2$ . □

For some reason, this proof is almost always presented as a proof by infinite descent. Notice that the induction hypothesis assumed that *both*  $p$  and  $q$  were as small as possible. Notice also that the ‘base cases’ included every pair of integers  $p$  and  $q$  where at least one of the integers is odd.

## 9 Fibonacci Parity

The *Fibonacci numbers*  $0, 1, 1, 2, 3, 5, 8, 13, 21, 34, 55, 89, 144, \dots$  are recursively defined as follows:

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

**Theorem 8.** For all non-negative integers  $n$ ,  $F_n$  is even if and only if  $n$  is divisible by 3.

**Proof:** Let  $n$  be an arbitrary non-negative integer. Assume that for all non-negative integers  $k < n$ ,  $F_k$  is even if and only if  $n$  is divisible by 3. There are three cases to consider:  $n = 0$ ,  $n = 1$ , and  $n \geq 2$ .

- If  $n = 0$ , then  $n$  is divisible by 3, and  $F_n = 0$  is even.
- If  $n = 1$ , then  $n$  is not divisible by 3, and  $F_n = 1$  is odd.
- If  $n \geq 2$ , there are two subcases to consider: Either  $n$  is divisible by 3, or it isn’t.
  - Suppose  $n$  is divisible by 3. Then neither  $n - 1$  nor  $n - 2$  is divisible by 3. Thus, the inductive hypothesis implies that both  $F_{n-1}$  and  $F_{n-2}$  are odd. So  $F_n$  is the sum of two odd numbers, and is therefore even.
  - Suppose  $n$  is not divisible by 3. Then exactly one of the numbers  $n - 1$  and  $n - 2$  is divisible by 3. Thus, the inductive hypothesis implies that exactly one of the numbers  $F_{n-1}$  and  $F_{n-2}$  is even, and the other is odd. So  $F_n$  is the sum of an even number and an odd number, and is therefore odd.

In all cases,  $F_n$  is even if and only if  $n$  is divisible by 3. □

## 10 Recursive Functions

**Theorem 9.** Suppose the function  $F: \mathbb{N} \rightarrow \mathbb{N}$  is defined recursively by setting  $F(0) = 0$  and  $F(n) = 1 + F(\lfloor n/2 \rfloor)$  for every positive integer  $n$ . Then for every positive integer  $n$ , we have  $F(n) = 1 + \lfloor \log_2 n \rfloor$ .

**Proof:** Let  $n$  be an arbitrary positive integer. Assume that  $F(k) = 1 + \lfloor \log_2 k \rfloor$  for every positive integer  $k < n$ . There are two cases to consider: Either  $n = 1$  or  $n \geq 2$ .

- Suppose  $n = 1$ . Then  $F(n) = F(1) = 1 + F(\lfloor 1/2 \rfloor) = 1 + F(0) = 1$  and  $1 + \lfloor \log_2 n \rfloor = 1 + \lfloor \log_2 1 \rfloor = 1 + \lfloor 0 \rfloor = 1$ .
- Suppose  $n \geq 2$ . Because  $1 \leq \lfloor n/2 \rfloor < n$ , the induction hypothesis implies that  $F(\lfloor n/2 \rfloor) = 1 + \lfloor \log_2 \lfloor n/2 \rfloor \rfloor$ . The definition of  $F(n)$  now implies that  $F(n) = 1 + F(\lfloor n/2 \rfloor) = 2 + \lfloor \log_2 \lfloor n/2 \rfloor \rfloor$ .

Now there are two subcases to consider:  $n$  is either even or odd.

- If  $n$  is even, then  $\lfloor n/2 \rfloor = n/2$ , which implies

$$\begin{aligned} F(n) &= 2 + \lfloor \log_2 \lfloor n/2 \rfloor \rfloor \\ &= 2 + \lfloor \log_2 (n/2) \rfloor \\ &= 2 + \lfloor (\log_2 n) - 1 \rfloor \\ &= 2 + \lfloor \log_2 n \rfloor - 1 \\ &= 1 + \lfloor \log_2 n \rfloor. \end{aligned}$$

- If  $n$  is odd, then  $\lfloor n/2 \rfloor = (n-1)/2$ , which implies

$$\begin{aligned} F(n) &= 2 + \lfloor \log_2 \lfloor n/2 \rfloor \rfloor \\ &= 2 + \lfloor \log_2 ((n-1)/2) \rfloor \\ &= 1 + \lfloor \log_2 (n-1) \rfloor \\ &= 1 + \lfloor \log_2 n \rfloor \end{aligned}$$

by the algebra in the even case. Because  $n > 1$  and  $n$  is odd,  $n$  cannot be a power of 2; thus,  $\lfloor \log_2 n \rfloor = \lfloor \log_2 (n-1) \rfloor$ .

In all cases, we conclude that  $F(n) = 1 + \lfloor \log_2 n \rfloor$ . □

## 11 Trees

Recall that a *tree* is a connected undirected graph with no cycles. A *subtree* of a tree  $T$  is a connected subgraph of  $T$ ; a *proper subtree* is any tree except  $T$  itself.

**Theorem 10.** In every tree, the number of vertices is one more than the number of edges.

This one is actually pretty easy to prove directly from the definition of ‘tree’: a connected acyclic graph.

**Proof:** Let  $T$  be an arbitrary tree. Choose an arbitrary vertex  $v$  of  $T$  to be the root, and direct every edge of  $T$  outward from  $v$ . Because  $T$  is connected, every node except  $v$  has at least one edge directed into it. Because  $T$  is acyclic, every node has at most one edge directed into it, and no edge is directed into  $v$ . Thus, for every node  $x \neq v$ , there is exactly one edge directed into  $x$ . We conclude that the number of edges is one less than the number of nodes. □

But we can prove this theorem by induction as well, in several different ways. Each inductive proof is structured around a different recursive definition of ‘tree’. First, a tree is either a single node, or two trees joined by an edge.

**Proof:** Let  $T$  be an arbitrary tree. Assume that in any proper subtree of  $T$ , the number of vertices is one more than the number of edges. There are two cases to consider: Either  $T$  has one vertex, or  $T$  has more than one vertex.

- If  $T$  has one vertex, then it has no edges.
- Suppose  $T$  has more than one vertex. Because  $T$  is connected, every pair of vertices is joined by a path. Thus,  $T$  must contain at least one edge. Let  $e$  be an arbitrary edge of  $T$ , and consider the graph  $T \setminus e$  obtained by deleting  $e$  from  $T$ .

Because  $T$  is acyclic, there is no path in  $T \setminus e$  between the endpoints of  $e$ . Thus,  $T$  has at least two connected components. On the other hand, because  $T$  is connected,  $T \setminus e$  has at most two connected components. Thus,  $T \setminus e$  has exactly two connected components; call them  $A$  and  $B$ .

Because  $T$  is acyclic, subgraphs  $A$  and  $B$  are also acyclic. Thus,  $A$  and  $B$  are subtrees of  $T$ , and therefore the induction hypothesis implies that  $|E(A)| = |V(A)| - 1$  and  $|E(B)| = |V(B)| - 1$ .

Because  $A$  and  $B$  do not share any vertices or edges, we have  $|V(T)| = |V(A)| + |V(B)|$  and  $|E(T)| = |E(A)| + |E(B)| + 1$ .

Simple algebra now implies that  $|E(T)| = |V(T)| - 1$ .

In both cases, we conclude that the number of vertices in  $T$  is one more than the number of edges in  $T$ .  $\square$

Second, a tree is a single node connected by edges to a finite set of trees.

**Proof:** Let  $T$  be an arbitrary tree. Assume that in any proper subtree of  $T$ , the number of vertices is one more than the number of edges. There are two cases to consider: Either  $T$  has one vertex, or  $T$  has more than one vertex.

- If  $T$  has one vertex, then it has no edges.
- Suppose  $T$  has more than one vertex. Let  $v$  be an arbitrary vertex of  $T$ , and let  $d$  be the degree of  $v$ . Delete  $v$  and all its incident edges from  $T$  to obtain a new graph  $G$ . This graph has exactly  $d$  connected components; call them  $G_1, G_2, \dots, G_d$ . Because  $T$  is acyclic, every subgraph of  $T$  is acyclic. Thus, every subgraph  $G_i$  is a proper subtree of  $G$ . So the induction hypothesis implies that  $|E(G_i)| = |V(G_i)| - 1$  for each  $i$ . We conclude that

$$|E(T)| = d + \sum_{i=1}^d |E(G_i)| = d + \sum_{i=1}^d (|V(G_i)| - 1) = \sum_{i=1}^d |V(G_i)| = |V(T)| - 1.$$

In both cases, we conclude that the number of vertices in  $T$  is one more than the number of edges in  $T$ .  $\square$

But you should **never** attempt to argue like this:

**Not a Proof:** The theorem is clearly true for the 1-node tree. So let  $T$  be an arbitrary tree with at least two nodes. Assume inductively that the number of vertices in  $T$  is one more than the number of edges in  $T$ . Suppose we add one more leaf to  $T$  to get a new tree  $T'$ . This new tree has one more vertex than  $T$  and one more edge than  $T$ . Thus, the number of vertices in  $T'$  is one more than the number of edges in  $T'$ .  $\square$

This is not a proof. Every sentence is true, and the connecting logic is correct, but it does not imply the theorem, because it doesn't *explicitly* consider *all possible* trees. Why should the reader believe that their favorite tree can be recursively constructed by adding leaves to a 1-node tree? It's *true*, of course, but that argument doesn't *prove* it. Remember: ***There are only two ways to prove any universally quantified statement:*** Directly ("Let  $T$  be an arbitrary tree...") or by contradiction ("Suppose some tree  $T$  doesn't...").

Here is a *correct* inductive proof using the same underlying idea. In this proof, I don't have to prove that the proof considers arbitrary trees; it says so right there on the first line! As usual, the proof very strongly resembles a recursive algorithm, including a subroutine to find a leaf.

**Proof:** Let  $T$  be an arbitrary tree. Assume that in any proper subtree of  $T$ , the number of vertices is one more than the number of edges. There are two cases to consider: Either  $T$  has one vertex, or  $T$  has more than one vertex.

- If  $T$  has one vertex, then it has no edges.
- Otherwise,  $T$  must have at least one vertex of degree 1, otherwise known as a leaf.

**Proof:** Consider a walk through the graph  $T$  that starts at an arbitrary vertex and continues as long as possible without repeating any edge. The walk can never visit the same vertex more than once, because  $T$  is acyclic. Whenever the walk visits a vertex of degree at least 2, it can continue further, because that vertex has at least one unvisited edge. But the walk must eventually end, because  $T$  is finite. Thus, the walk must eventually reach a vertex of degree 1.  $\square$

Let  $\ell$  be an arbitrary leaf of  $T$ , and let  $T'$  be the tree obtained by deleting  $\ell$  from  $T$ . Then we have the identity

$$|E(T)| = |E(T')| - 1 = |V(T')| - 2 = |V(T)| - 1,$$

where the first and third equalities follow from the definition of  $T'$ , and the second equality follows from the inductive hypothesis.

In both cases, we conclude that the number of vertices in  $T$  is one more than the number of edges in  $T$ .  $\square$

## 12 Strings

Recall that a *string* is any finite sequence of symbols. More formally, a string is either empty or a single symbol followed by a string. Let  $x \cdot y$  denote the concatenation of strings  $x$  and  $y$ , and let  $\text{reverse}(x)$  denote the reversal of string  $x$ . For example,  $\text{now} \cdot \text{here} = \text{nowhere}$  and  $\text{reverse}(\text{stop}) = \text{pots}$ .

**Theorem 11.** For all strings  $x$  and  $y$ , we have  $\text{reverse}(x \cdot y) = \text{reverse}(y) \cdot \text{reverse}(x)$ .

**Proof:** Let  $x$  and  $y$  be arbitrary strings. Assume for any proper substring  $w$  of  $x$  that  $\text{reverse}(w \cdot y) = \text{reverse}(y) \cdot \text{reverse}(w)$ . There are two cases to consider: Either  $x$  is the empty string, or not.

- If  $x$  is the empty string, then  $\text{reverse}(x)$  is also the empty string, so  $\text{reverse}(x \cdot y) = \text{reverse}(y) = \text{reverse}(y) \cdot \text{reverse}(x)$ .
- Suppose  $x$  is not the empty string. Then  $x = a \cdot w$ , for some character  $a$  and some string  $w$ , and therefore

$$\begin{aligned}\text{reverse}(x \cdot y) &= \text{reverse}((a \cdot w) \cdot y) \\ &= \text{reverse}(a \cdot (w \cdot y)) \\ &= \text{reverse}(w \cdot y) \cdot a \\ &= \text{reverse}(y) \cdot \text{reverse}(w) \cdot a \quad [\text{induction hypothesis}] \\ &= \text{reverse}(y) \cdot \text{reverse}(x).\end{aligned}$$

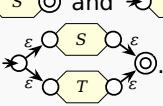
In both cases, we conclude that  $\text{reverse}(x \cdot y) = \text{reverse}(y) \cdot \text{reverse}(x)$ .  $\square$

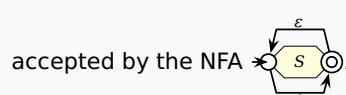
## 13 Regular Languages

**Theorem 12.** Every regular language is accepted by a non-deterministic finite automaton.

**Proof:** In fact, we will show something stronger: Every regular language is accepted by an NFA with exactly one accepting state.

Let  $R$  be an arbitrary regular expression over the finite alphabet  $\Sigma$ . Assume that for any sub-expression  $S$  of  $R$ , the corresponding regular language is accepted by an NFA with one accepting state, denoted . There are six cases to consider—three base cases and three recursive cases—mirroring the recursive definition of a regular expression.

- If  $R = \emptyset$ , then  $L(R) = \emptyset$  is accepted by an NFA with no transitions: .
- If  $R = \epsilon$ , then  $L(R) = \{\epsilon\}$  is accepted by the NFA .
- If  $R = a$  for some character  $a \in \Sigma$ , then  $L(R) = \{a\}$  is accepted by the NFA .
- Suppose  $R = ST$  for some regular expressions  $S$  and  $T$ . The inductive hypothesis implies that  $S$  and  $T$  are accepted by NFAs  and , respectively. Then  $L(T) = \{uv \mid u \in L(S), v \in L(T)\}$  is accepted by the NFA .
- Suppose  $R = S + T$  for some regular expressions  $S$  and  $T$ . The inductive hypothesis implies that  $S$  and  $T$  are accepted by NFAs  and , respectively. Then  $L(R) = L(S) \cup L(T)$  is accepted by the NFA .
- Finally, suppose  $R = S^*$  for some regular expression  $S$ . The inductive hypothesis implies that  $S$  is accepted by an NFA .



In every case,  $L(x)$  is accepted by an NFA with one accepting state.  $\square$

*Change is certain. Peace is followed by disturbances; departure of evil men by their return. Such recurrences should not constitute occasions for sadness but realities for awareness, so that one may be happy in the interim.*

— I Ching [The Book of Changes] (c. 1100 BC)

*“... O Zarathustra, who you are and must become” behold you are the teacher of the eternal recurrence – that is your destiny! That you as the first must teach this doctrine – how could this great destiny not be your greatest danger and sickness too?*

— Friedrich Nietzsche, *Also sprach Zarathustra* (1885)  
[translated by Walter Kaufmann]

**Wil Wheaton:** Embrace the dark side!

**Sheldon:** That's not even from your franchise!

— “The Wheaton Recurrence”, *Bing Bang Theory*, April 12, 2010

# Solving Recurrences

## 1 Introduction

A **recurrence** is a recursive description of a function, or in other words, a description of a function in terms of itself. Like all recursive structures, a recurrence consists of one or more *base cases* and one or more *recursive cases*. Each of these cases is an equation or inequality, with some function value  $f(n)$  on the left side. The base cases give explicit values for a (typically finite, typically small) subset of the possible values of  $n$ . The recursive cases relate the function value  $f(n)$  to function value  $f(k)$  for one or more integers  $k < n$ ; typically, each recursive case applies to an infinite number of possible values of  $n$ .

For example, the following recurrence (written in two different but standard ways) describes the identity function  $f(n) = n$ :

$$f(n) = \begin{cases} 0 & \text{if } n = 0 \\ f(n-1) + 1 & \text{otherwise} \end{cases} \quad \begin{aligned} f(0) &= 0 \\ f(n) &= f(n-1) + 1 \text{ for all } n > 0 \end{aligned}$$

In both presentations, the first line is the only base case, and the second line is the only recursive case. The same function can satisfy *many* different recurrences; for example, both of the following recurrences also describe the identity function:

$$f(n) = \begin{cases} 0 & \text{if } n = 0 \\ 1 & \text{if } n = 1 \\ f(\lfloor n/2 \rfloor) + f(\lceil n/2 \rceil) & \text{otherwise} \end{cases} \quad f(n) = \begin{cases} 0 & \text{if } n = 0 \\ 2 \cdot f(n/2) & \text{if } n \text{ is even and } n > 0 \\ f(n-1) + 1 & \text{if } n \text{ is odd} \end{cases}$$

We say that a particular function **satisfies** a recurrence, or is the **solution** to a recurrence, if each of the statements in the recurrence is true. Most recurrences—at least, those that we will encounter in this class—have a solution; moreover, if every case of the recurrence is an equation, that solution is unique. Specifically, if we transform the recursive formula into a recursive *algorithm*, the solution to the recurrence is the function computed by that algorithm!

Recurrences arise naturally in the analysis of algorithms, especially recursive algorithms. In many cases, we can express the running time of an algorithm as a recurrence, where the recursive cases of the recurrence correspond exactly to the recursive cases of the algorithm. Recurrences are also useful tools for solving counting problems—How many objects of a particular kind exist?

By itself, a recurrence is not a satisfying description of the running time of an algorithm or a bound on the number of widgets. Instead, we need a *closed-form* solution to the recurrence; this is a *non-recursive* description of a function that satisfies the recurrence. For recurrence *equations*, we sometimes prefer an *exact* closed-form solution, but such a solution may not exist, or may be too complex to be useful. Thus, for most recurrences, especially those arising in algorithm analysis, we can be satisfied with an *asymptotic* solution of the form  $\Theta(f(n))$ , for some explicit (non-recursive) function  $g(n)$ .

For recursive *inequalities*, we prefer a *tight* solution; this is a function that would still satisfy the recurrence if all the inequalities were replaced with the corresponding equations. Again, exactly tight solutions may not exist, or may be too complex to be useful, so we may have to settle for a looser solution and/or an asymptotic solution of the form  $O(g(n))$  or  $\Omega(g(n))$ .

## 2 The Ultimate Method: Guess and Confirm

Ultimately, there is only one fail-safe method to solve *any* recurrence:

**Guess the answer, and then prove it correct by induction.**

Later sections of these notes describe techniques to generate guesses that are guaranteed to be correct, provided you use them correctly. But if you're faced with a recurrence that doesn't seem to fit any of these methods, or if you've forgotten how those techniques work, don't despair! If you guess a closed-form solution and then try to verify your guess inductively, usually either the proof will succeed, in which case you're done, or the proof will fail, in which case *the failure will help you refine your guess*. Where you get your initial guess is utterly irrelevant<sup>1</sup>—from a classmate, from a textbook, on the web, from the answer to a different problem, scrawled on a bathroom wall in Siebel, included in a care package from your mom, dictated by the machine elves, whatever. If you can prove that the answer is correct, then it's correct!

### 2.1 Tower of Hanoi

The classical Tower of Hanoi problem gives us the recurrence  $T(n) = 2T(n - 1) + 1$  with base case  $T(0) = 0$ . Just looking at the recurrence we can guess that  $T(n)$  is something like  $2^n$ . If we write out the first few values of  $T(n)$ , we discover that they are each one less than a power of two.

$$T(0) = 0, \quad T(1) = 1, \quad T(2) = 3, \quad T(3) = 7, \quad T(4) = 15, \quad T(5) = 31, \quad T(6) = 63, \quad \dots,$$

It looks like  $T(n) = 2^n - 1$  might be the right answer. Let's check.

$$\begin{aligned} T(0) &= 0 = 2^0 - 1 & \checkmark \\ T(n) &= 2T(n - 1) + 1 \\ &= 2(2^{n-1} - 1) + 1 & [\text{induction hypothesis}] \\ &= 2^n - 1 & [\text{algebra}] \end{aligned}$$

We were right! Hooray, we're done!

---

<sup>1</sup>... except of course during exams, where you aren't supposed to use *any* outside sources

Another way we can guess the solution is by *unrolling* the recurrence, by substituting it into itself:

$$\begin{aligned} T(n) &= 2T(n-1) + 1 \\ &= 2(2T(n-2) + 1) + 1 \\ &= 4T(n-2) + 3 \\ &= 4(2T(n-3) + 1) + 3 \\ &= 8T(n-3) + 7 \\ &= \dots \end{aligned}$$

It looks like unrolling the initial Hanoi recurrence  $k$  times, for any non-negative integer  $k$ , will give us the new recurrence  $T(n) = 2^k T(n-k) + (2^k - 1)$ . Let's prove this by induction:

$$\begin{aligned} T(n) &= 2T(n-1) + 1 \quad \checkmark & [k=0, \text{ by definition}] \\ T(n) &= 2^{k-1}T(n-(k-1)) + (2^{k-1} - 1) & [\text{inductive hypothesis}] \\ &= 2^{k-1}(2T(n-k) + 1) + (2^{k-1} - 1) & [\text{initial recurrence for } T(n-(k-1))] \\ &= 2^k T(n-k) + (2^k - 1) \quad \checkmark & [\text{algebra}] \end{aligned}$$

Our guess was correct! In particular, unrolling the recurrence  $n$  times give us the recurrence  $T(n) = 2^n T(0) + (2^n - 1)$ . Plugging in the base case  $T(0) = 0$  give us the closed-form solution  $T(n) = 2^n - 1$ .

## 2.2 Fibonacci numbers

Let's try a less trivial example: the Fibonacci numbers  $F_n = F_{n-1} + F_{n-2}$  with base cases  $F_0 = 0$  and  $F_1 = 1$ . There is no obvious pattern in the first several values (aside from the recurrence itself), but we can reasonably guess that  $F_n$  is exponential in  $n$ . Let's try to prove inductively that  $F_n \leq \alpha \cdot c^n$  for some constants  $\alpha > 0$  and  $c > 1$  and see how far we get.

$$\begin{aligned} F_n &= F_{n-1} + F_{n-2} \\ &\leq \alpha \cdot c^{n-1} + \alpha \cdot c^{n-2} & [\text{"induction hypothesis"}] \\ &\leq \alpha \cdot c^n \quad ??? \end{aligned}$$

The last inequality is satisfied if  $c^n \geq c^{n-1} + c^{n-2}$ , or more simply, if  $c^2 - c - 1 \geq 0$ . The smallest value of  $c$  that works is  $\phi = (1 + \sqrt{5})/2 \approx 1.618034$ ; the other root of the quadratic equation has smaller absolute value, so we can ignore it.

So we have *most* of an inductive proof that  $F_n \leq \alpha \cdot \phi^n$  for *some* constant  $\alpha$ . All that we're missing are the base cases, which (we can easily guess) must determine the value of the coefficient  $\alpha$ . We quickly compute

$$\frac{F_0}{\phi^0} = \frac{0}{1} = 0 \quad \text{and} \quad \frac{F_1}{\phi^1} = \frac{1}{\phi} \approx 0.618034 > 0,$$

so the base cases of our induction proof are correct as long as  $\alpha \geq 1/\phi$ . It follows that  $F_n \leq \phi^{n-1}$  for all  $n \geq 0$ .

What about a matching lower bound? Essentially the same inductive proof implies that  $F_n \geq \beta \cdot \phi^n$  for some constant  $\beta$ , but the only value of  $\beta$  that works for *all*  $n$  is the trivial  $\beta = 0$ ! We could try to find some lower-order term that makes the base case non-trivial, but an easier approach is to recall that

asymptotic  $\Omega()$  bounds only have to work for *sufficiently large*  $n$ . So let's ignore the trivial base case  $F_0 = 0$  and assume that  $F_2 = 1$  is a base case instead. Some more easy calculation gives us

$$\frac{F_2}{\phi^2} = \frac{1}{\phi^2} \approx 0.381966 < \frac{1}{\phi}.$$

Thus, the new base cases of our induction proof are correct as long as  $\beta \leq 1/\phi^2$ , which implies that  $F_n \geq \phi^{n-2}$  for all  $n \geq 1$ .

Putting the upper and lower bounds together, we obtain the tight asymptotic bound  $F_n = \Theta(\phi^n)$ . It is possible to get a more exact solution by speculatively refining and conforming our current bounds, but it's not easy. Fortunately, if we really need it, we can get an exact solution using the *annihilator* method, which we'll see later in these notes.

### 2.3 Mergesort

Mergesort is a classical recursive divide-and-conquer algorithm for sorting an array. The algorithm splits the array in half, recursively sorts the two halves, and then merges the two sorted subarrays into the final sorted array.

```
MERGESORT(A[1..n]):
  if (n > 1)
    m ← ⌊n/2⌋
    MERGESORT(A[1..m])
    MERGESORT(A[m+1..n])
    MERGE(A[1..n], m)
```

```
MERGE(A[1..n], m):
  i ← 1; j ← m + 1
  for k ← 1 to n
    if j > n
      B[k] ← A[i]; i ← i + 1
    else if i > m
      B[k] ← A[j]; j ← j + 1
    else if A[i] < A[j]
      B[k] ← A[i]; i ← i + 1
    else
      B[k] ← A[j]; j ← j + 1
  for k ← 1 to n
    A[k] ← B[k]
```

Let  $T(n)$  denote the worst-case running time of MERGESORT when the input array has size  $n$ . The MERGE subroutine clearly runs in  $\Theta(n)$  time, so the function  $T(n)$  satisfies the following recurrence:

$$T(n) = \begin{cases} \Theta(1) & \text{if } n = 1, \\ T(\lceil n/2 \rceil) + T(\lfloor n/2 \rfloor) + \Theta(n) & \text{otherwise.} \end{cases}$$

For now, let's consider the special case where  $n$  is a power of 2; this assumption allows us to take the floors and ceilings out of the recurrence. (We'll see how to deal with the floors and ceilings later; the short version is that they don't matter.)

Because the recurrence itself is given only asymptotically—in terms of  $\Theta()$  expressions—we can't hope for anything but an asymptotic solution. So we can safely simplify the recurrence further by removing the  $\Theta$ 's; any asymptotic solution to the simplified recurrence will also satisfy the original recurrence. (This simplification is actually important for another reason; if we kept the asymptotic expressions, we might be tempted to simplify them inappropriately.)

Our simplified recurrence now looks like this:

$$T(n) = \begin{cases} 1 & \text{if } n = 1, \\ 2T(n/2) + n & \text{otherwise.} \end{cases}$$

To guess at a solution, let's try unrolling the recurrence.

$$\begin{aligned} T(n) &= 2T(n/2) + n \\ &= 2(2T(n/4) + n/2) + n \\ &= 4T(n/4) + 2n \\ &= 8T(n/8) + 3n = \dots \end{aligned}$$

It looks like  $T(n)$  satisfies the recurrence  $T(n) = 2^k T(n/2^k) + kn$  for any positive integer  $k$ . Let's verify this by induction.

$$\begin{aligned} T(n) &= 2T(n/2) + n = 2^1 T(n/2^1) + 1 \cdot n \quad \checkmark & [k = 1, \text{ given recurrence}] \\ T(n) &= 2^{k-1} T(n/2^{k-1}) + (k-1)n & [\text{inductive hypothesis}] \\ &= 2^{k-1} (2T(n/2^k) + n/2^{k-1}) + (k-1)n & [\text{substitution}] \\ &= 2^k T(n/2^k) + kn \quad \checkmark & [\text{algebra}] \end{aligned}$$

Our guess was right! The recurrence becomes trivial when  $n/2^k = 1$ , or equivalently, when  $k = \log_2 n$ :

$$T(n) = nT(1) + n \log_2 n = n \log_2 n + n.$$

Finally, we have to put back the  $\Theta$ 's we stripped off; our final closed-form solution is  $T(n) = \Theta(n \log n)$ .

## 2.4 An uglier divide-and-conquer example

Consider the divide-and-conquer recurrence  $T(n) = \sqrt{n} \cdot T(\sqrt{n}) + n$ . This doesn't fit into the form required by the Master Theorem (which we'll see below), but it still sort of resembles the Mergesort recurrence—the total size of the subproblems at the first level of recursion is  $n$ —so let's guess that  $T(n) = O(n \log n)$ , and then try to prove that our guess is correct. (We could also attack this recurrence by unrolling, but let's see how far just guessing will take us.)

Let's start by trying to prove an upper bound  $T(n) \leq a n \lg n$  for all sufficiently large  $n$  and some constant  $a$  to be determined later:

$$\begin{aligned} T(n) &= \sqrt{n} \cdot T(\sqrt{n}) + n \\ &\leq \sqrt{n} \cdot a \sqrt{n} \lg \sqrt{n} + n & [\text{induction hypothesis}] \\ &= (a/2)n \lg n + n & [\text{algebra}] \\ &\leq an \lg n \quad \checkmark & [\text{algebra}] \end{aligned}$$

The last inequality assumes only that  $1 \leq (a/2) \log n$ , or equivalently, that  $n \geq 2^{2/a}$ . In other words, the induction proof is correct if  $n$  is sufficiently large. So we were right!

But before you break out the champagne, what about the multiplicative constant  $a$ ? The proof worked for any constant  $a$ , no matter how small. This strongly suggests that our upper bound  $T(n) = O(n \log n)$  is not tight. Indeed, if we try to prove a matching lower bound  $T(n) \geq b n \log n$  for sufficiently large  $n$ , we run into trouble.

$$\begin{aligned} T(n) &= \sqrt{n} \cdot T(\sqrt{n}) + n \\ &\geq \sqrt{n} \cdot b \sqrt{n} \log \sqrt{n} + n & [\text{induction hypothesis}] \\ &= (b/2)n \log n + n \\ &\not\geq bn \log n \end{aligned}$$

The last inequality would be correct only if  $1 > (b/2) \log n$ , but that inequality is false for large values of  $n$ , no matter which constant  $b$  we choose.

Okay, so  $\Theta(n \log n)$  is too big. How about  $\Theta(n)$ ? The lower bound is easy to prove directly:

$$T(n) = \sqrt{n} \cdot T(\sqrt{n}) + n \geq n \checkmark$$

But an inductive proof of the upper bound fails.

$$\begin{aligned} T(n) &= \sqrt{n} \cdot T(\sqrt{n}) + n \\ &\leq \sqrt{n} \cdot a \sqrt{n} + n && [\text{induction hypothesis}] \\ &= (a+1)n && [\text{algebra}] \\ &\not\leq an \end{aligned}$$

Hmmm. So what's bigger than  $n$  and smaller than  $n \lg n$ ? How about  $n \sqrt{\lg n}$ ?

$$\begin{aligned} T(n) &= \sqrt{n} \cdot T(\sqrt{n}) + n \leq \sqrt{n} \cdot a \sqrt{n} \sqrt{\lg \sqrt{n}} + n && [\text{induction hypothesis}] \\ &= (a/\sqrt{2}) n \sqrt{\lg n} + n && [\text{algebra}] \\ &\leq a n \sqrt{\lg n} \quad \text{for large enough } n \checkmark \end{aligned}$$

Okay, the upper bound checks out; how about the lower bound?

$$\begin{aligned} T(n) &= \sqrt{n} \cdot T(\sqrt{n}) + n \geq \sqrt{n} \cdot b \sqrt{n} \sqrt{\lg \sqrt{n}} + n && [\text{induction hypothesis}] \\ &= (b/\sqrt{2}) n \sqrt{\lg n} + n && [\text{algebra}] \\ &\not\geq b n \sqrt{\lg n} \end{aligned}$$

No, the last step doesn't work. So  $\Theta(n \sqrt{\lg n})$  doesn't work.

Okay... what else is between  $n$  and  $n \lg n$ ? How about  $n \lg \lg n$ ?

$$\begin{aligned} T(n) &= \sqrt{n} \cdot T(\sqrt{n}) + n \leq \sqrt{n} \cdot a \sqrt{n} \lg \lg \sqrt{n} + n && [\text{induction hypothesis}] \\ &= a n \lg \lg n - a n + n && [\text{algebra}] \\ &\leq a n \lg \lg n \quad \text{if } a \geq 1 \checkmark \end{aligned}$$

Hey look at that! For once, our upper bound proof requires a constraint on the hidden constant  $a$ . This is a good indication that we've found the right answer. Let's try the lower bound:

$$\begin{aligned} T(n) &= \sqrt{n} \cdot T(\sqrt{n}) + n \geq \sqrt{n} \cdot b \sqrt{n} \lg \lg \sqrt{n} + n && [\text{induction hypothesis}] \\ &= b n \lg \lg n - b n + n && [\text{algebra}] \\ &\geq b n \lg \lg n \quad \text{if } b \leq 1 \checkmark \end{aligned}$$

Hey, it worked! We have most of an inductive proof that  $T(n) \leq a n \lg \lg n$  for any  $a \geq 1$  and most of an inductive proof that  $T(n) \geq b n \lg \lg n$  for any  $b \leq 1$ . Technically, we're still missing the base cases in both proofs, but we can be fairly confident at this point that  $T(n) = \Theta(n \log \log n)$ .

### 3 Divide and Conquer Recurrences (Recursion Trees)

Many divide and conquer algorithms give us running-time recurrences of the form

$$T(n) = a T(n/b) + f(n) \quad (1)$$

where  $a$  and  $b$  are constants and  $f(n)$  is some other function. There is a simple and general technique for solving many recurrences in this and similar forms, using a **recursion tree**. The root of the recursion tree is a box containing the value  $f(n)$ ; the root has  $a$  children, each of which is the root of a (recursively defined) recursion tree for the function  $T(n/b)$ .

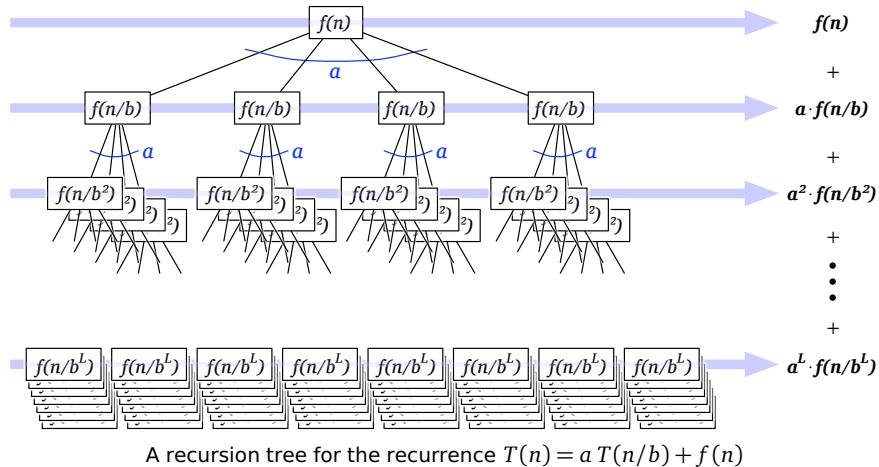
Equivalently, a recursion tree is a complete  $a$ -ary tree where each node at depth  $i$  contains the value  $f(n/b^i)$ . The recursion stops when we get to the base case(s) of the recurrence. Because we're only looking for asymptotic bounds, the exact base case doesn't matter; we can safely assume that  $T(1) = \Theta(1)$ , or even that  $T(n) = \Theta(1)$  for all  $n \leq 10^{100}$ . I'll also assume for simplicity that  $n$  is an integral power of  $b$ ; we'll see how to avoid this assumption later (but to summarize: it doesn't matter).

Now  $T(n)$  is just the sum of all values stored in the recursion tree. For each  $i$ , the  $i$ th level of the tree contains  $a^i$  nodes, each with value  $f(n/b^i)$ . Thus,

$$T(n) = \sum_{i=0}^L a^i f(n/b^i) \quad (\Sigma)$$

where  $L$  is the depth of the recursion tree. We easily see that  $L = \log_b n$ , because  $n/b^L = 1$ . The base case  $f(1) = \Theta(1)$  implies that the last non-zero term in the summation is  $\Theta(a^L) = \Theta(a^{\log_b n}) = \Theta(n^{\log_b a})$ .

For most divide-and-conquer recurrences, the level-by-level sum (??) is a geometric series—each term is a constant factor larger or smaller than the previous term. In this case, only the largest term in the geometric series matters; all of the other terms are swallowed up by the  $\Theta(\cdot)$  notation.



Here are several examples of the recursion-tree technique in action:

- **Mergesort (simplified):**  $T(n) = 2T(n/2) + n$

There are  $2^i$  nodes at level  $i$ , each with value  $n/2^i$ , so every term in the level-by-level sum (??) is the same:

$$T(n) = \sum_{i=0}^L n.$$

The recursion tree has  $L = \log_2 n$  levels, so  $T(n) = \Theta(n \log n)$ .

- **Randomized selection:**  $T(n) = T(3n/4) + n$

The recursion tree is a single path. The node at depth  $i$  has value  $(3/4)^i n$ , so the level-by-level sum (??) is a decreasing geometric series:

$$T(n) = \sum_{i=0}^L (3/4)^i n.$$

This geometric series is dominated by its initial term  $n$ , so  $T(n) = \Theta(n)$ . The recursion tree has  $L = \log_{4/3} n$  levels, but so what?

- **Karatsuba's multiplication algorithm:**  $T(n) = 3T(n/2) + n$

There are  $3^i$  nodes at depth  $i$ , each with value  $n/2^i$ , so the level-by-level sum (??) is an increasing geometric series:

$$T(n) = \sum_{i=0}^L (3/2)^i n.$$

This geometric series is dominated by its final term  $(3/2)^L n$ . Each leaf contributes 1 to this term; thus, the final term is equal to the number of leaves in the tree! The recursion tree has  $L = \log_2 n$  levels, and therefore  $3^{\log_2 n} = n^{\log_2 3}$  leaves, so  $T(n) = \Theta(n^{\log_2 3})$ .

- $T(n) = 2T(n/2) + n / \lg n$

The sum of all the nodes in the  $i$ th level is  $n/(\lg n - i)$ . This implies that the depth of the tree is at most  $\lg n - 1$ . The level sums are neither constant nor a geometric series, so we just have to evaluate the overall sum directly.

Recall (or if you're seeing this for the first time: Behold!) that the  $n$ th *harmonic number*  $H_n$  is the sum of the reciprocals of the first  $n$  positive integers:

$$H_n := \sum_{i=1}^n \frac{1}{i}$$

It's not hard to show that  $H_n = \Theta(\log n)$ ; in fact, we have the stronger inequalities  $\ln(n+1) \leq H_n \leq \ln n + 1$ .

$$T(n) = \sum_{i=0}^{\lg n - 1} \frac{n}{\lg n - i} = \sum_{j=1}^{\lg n} \frac{n}{j} = nH_{\lg n} = \Theta(n \lg \lg n)$$

- $T(n) = 4T(n/2) + n \lg n$

There are  $4^i$  nodes at each level  $i$ , each with value  $(n/2^i) \lg(n/2^i) = (n/2^i)(\lg n - i)$ ; again, the depth of the tree is at most  $\lg n - 1$ . We have the following summation:

$$T(n) = \sum_{i=0}^{\lg n - 1} n2^i (\lg n - i)$$

We can simplify this sum by substituting  $j = \lg n - i$ :

$$T(n) = \sum_{j=i}^{\lg n} n2^{\lg n - j} j = \sum_{j=i}^{\lg n} n^2 j / 2^j = \Theta(j^2)$$

The last step uses the fact that  $\sum_{i=1}^{\infty} j/2^j = 2$ . Although this is not quite a geometric series, it is still dominated by its largest term.

- **Ugly divide and conquer:**  $T(n) = \sqrt{n} \cdot T(\sqrt{n}) + n$

We solved this recurrence earlier by guessing the right answer and verifying, but we can use recursion trees to get the correct answer directly. The *degree* of the nodes in the recursion tree is no longer constant, so we have to be a bit more careful, but the same basic technique still applies. It's not hard to see that the nodes in any level sum to  $n$ . The depth  $L$  satisfies the identity  $n^{2^{-L}} = 2$  (we can't get all the way down to 1 by taking square roots), so  $L = \lg \lg n$  and  $T(n) = \Theta(n \lg \lg n)$ .

- **Randomized quicksort:**  $T(n) = T(3n/4) + T(n/4) + n$

This recurrence isn't in the standard form described earlier, but we can still solve it using recursion trees. Now nodes in the same level of the recursion tree have different values, and different leaves are at different levels. However, the nodes in any *complete* level (that is, above any of the leaves) sum to  $n$ . Moreover, every leaf in the recursion tree has depth between  $\log_4 n$  and  $\log_{4/3} n$ . To derive an upper bound, we overestimate  $T(n)$  by ignoring the base cases and extending the tree downward to the level of the *deepest* leaf. Similarly, to derive a lower bound, we overestimate  $T(n)$  by counting only nodes in the tree up to the level of the *shallowest* leaf. These observations give us the upper and lower bounds  $n \log_4 n \leq T(n) \leq n \log_{4/3} n$ . Since these bounds differ by only a constant factor, we have  $T(n) = \Theta(n \log n)$ .

- **Deterministic selection:**  $T(n) = T(n/5) + T(7n/10) + n$

Again, we have a lopsided recursion tree. If we look only at complete levels of the tree, we find that the level sums form a descending geometric series  $T(n) = n + 9n/10 + 81n/100 + \dots$ . We can get an upper bound by ignoring the base cases entirely and growing the tree out to infinity, and we can get a lower bound by only counting nodes in complete levels. Either way, the geometric series is dominated by its largest term, so  $T(n) = \Theta(n)$ .

- **Randomized search trees:**  $T(n) = \frac{1}{4}T(n/4) + \frac{3}{4}T(3n/4) + 1$

This looks like a divide-and-conquer recurrence, but what does it mean to have a quarter of a child? The right approach is to imagine that each node in the recursion tree has a *weight* in addition to its value. Alternately, we get a standard recursion tree again if we add a second real parameter to the recurrence, defining  $T(n) = T(n, 1)$ , where

$$T(n, \alpha) = T(n/4, \alpha/4) + T(3n/4, 3\alpha/4) + \alpha.$$

In each complete level of the tree, the (weighted) node values sum to exactly 1. The leaves of the recursion tree are at different levels, but all between  $\log_4 n$  and  $\log_{4/3} n$ . So we have upper and lower bounds  $\log_4 n \leq T(n) \leq \log_{4/3} n$ , which differ by only a constant factor, so  $T(n) = \Theta(\log n)$ .

- **Ham-sandwich trees:**  $T(n) = T(n/2) + T(n/4) + 1$

Again, we have a lopsided recursion tree. If we only look at complete levels, we find that the level sums form an *ascending* geometric series  $T(n) = 1 + 2 + 4 + \dots$ , so the solution is dominated by the number of leaves. The recursion tree has  $\log_4 n$  complete levels, so there are more than

$2^{\log_4 n} = n^{\log_4 2} = \sqrt{n}$ ; on the other hand, every leaf has depth at most  $\log_2 n$ , so the total number of leaves is at most  $2^{\log_2 n} = n$ . Unfortunately, the crude bounds  $\sqrt{n} \ll T(n) \ll n$  are the best we can derive using the techniques we know so far!

The following theorem completely describes the solution for any divide-and-conquer recurrence in the ‘standard form’  $T(n) = aT(n/b) + f(n)$ , where  $a$  and  $b$  are constants and  $f(n)$  is a polynomial. This theorem allows us to bypass recursion trees for ‘standard’ recurrences, but many people (including Jeff) find it harder to remember than the more general recursion-tree technique. Your mileage may vary.

**The Master Theorem.** *The recurrence  $T(n) = aT(n/b) + f(n)$  can be solved as follows.*

- If  $a f(n/b) = \kappa f(n)$  for some constant  $\kappa < 1$ , then  $T(n) = \Theta(f(n))$ .
- If  $a f(n/b) = K f(n)$  for some constant  $K > 1$ , then  $T(n) = \Theta(n^{\log_b a})$ .
- If  $a f(n/b) = f(n)$ , then  $T(n) = \Theta(f(n) \log_b n)$ .
- If none of these three cases apply, you’re on your own.

**Proof:** If  $f(n)$  is a *constant factor larger* than  $a f(b/n)$ , then by induction, the sum is a descending geometric series. The sum of any geometric series is a constant times its largest term. In this case, the largest term is the first term  $f(n)$ .

If  $f(n)$  is a *constant factor smaller* than  $a f(b/n)$ , then by induction, the sum is an ascending geometric series. The sum of any geometric series is a constant times its largest term. In this case, this is the last term, which by our earlier argument is  $\Theta(n^{\log_b a})$ .

Finally, if  $a f(b/n) = f(n)$ , then by induction, each of the  $L+1$  terms in the sum is equal to  $f(n)$ .  $\square$

## \*4 The Nuclear Bomb

Finally, let me describe *without proof* a powerful generalization of the recursion tree method, first published by Lebanese researchers Mohamad Akra and Louay Bazzi in 1998. Consider a general divide-and-conquer recurrence of the form

$$T(n) = \sum_{i=1}^k a_i T(n/b_i) + f(n),$$

where  $k$  is a constant,  $a_i > 0$  and  $b_i > 1$  are constants for all  $i$ , and  $f(n) = \Omega(n^c)$  and  $f(n) = O(n^d)$  for some constants  $0 < c \leq d$ . (As usual, we assume the standard base case  $T(\Theta(1)) = \Theta(1)$ .) Akra and Bazzi prove that this recurrence has the closed-form asymptotic solution

$$T(n) = \Theta\left(n^\rho \left(1 + \int_1^n \frac{f(u)}{u^{\rho+1}} du\right)\right),$$

where  $\rho$  is the unique real solution to the equation

$$\sum_{i=1}^k a_i / b_i^\rho = 1.$$

In particular, the Akra-Bazzi theorem immediately implies the following form of the Master Theorem:

$$T(n) = aT(n/b) + n^c \implies T(n) = \begin{cases} \Theta(n^{\log_b a}) & \text{if } c < \log_b a - \varepsilon \\ \Theta(n^c \log n) & \text{if } c = \log_b a \\ \Theta(n^c) & \text{if } c > \log_b a + \varepsilon \end{cases}$$

The Akra-Bazzi theorem does not require that the parameters  $a_i$  and  $b_i$  are integers, or even rationals; on the other hand, even when all parameters are integers, the characteristic equation  $\sum_i a_i/b_i^\rho = 1$  may have no analytical solution.

Here are a few examples of recurrences that are difficult (or impossible) for recursion trees, but have easy solutions using the Akra-Bazzi theorem.

- **Randomized quicksort:**  $T(n) = T(3n/4) + T(n/4) + n$

The equation  $(3/4)^\rho + (1/4)^\rho = 1$  has the unique solution  $\rho = 1$ , and therefore

$$T(n) = \Theta\left(n\left(1 + \int_1^n \frac{1}{u} du\right)\right) = O(n \log n).$$

- **Deterministic selection:**  $T(n) = T(n/5) + T(7n/10) + n$

The equation  $(1/5)^\rho + (7/10)^\rho = 1$  has no analytical solution. However, we easily observe that  $(1/5)^x + (7/10)^x$  is a decreasing function of  $x$ , and therefore  $0 < \rho < 1$ . Thus, we have

$$\int_1^n \frac{f(u)}{u^{\rho+1}} du = \int_1^n u^{-\rho} du = \frac{u^{1-\rho}}{1-\rho} \Big|_{u=1}^n = \frac{n^{1-\rho} - 1}{1-\rho} = \Theta(n^{1-\rho}),$$

and therefore

$$T(n) = \Theta(n^\rho \cdot (1 + \Theta(n^{1-\rho}))) = \Theta(n).$$

- **Randomized search trees:**  $T(n) = \frac{1}{4}T(n/4) + \frac{3}{4}T(3n/4) + 1$

The equation  $\frac{1}{4}(\frac{1}{4})^\rho + \frac{3}{4}(\frac{3}{4})^\rho = 1$  has the unique solution  $\rho = 0$ , and therefore

$$T(n) = \Theta\left(1 + \int_1^n \frac{1}{u} du\right) = \Theta(\log n).$$

- **Ham-sandwich trees:**  $T(n) = T(n/2) + T(n/4) + 1$ . Recall that we could only prove the very weak bounds  $\sqrt{n} \ll T(n) \ll n$  using recursion trees. The equation  $(1/2)^\rho + (1/4)^\rho = 1$  has the unique solution  $\rho = \log_2((1 + \sqrt{5})/2) \approx 0.69424$ , which can be obtained by setting  $x = 2^\rho$  and solving for  $x$ . Thus, we have

$$\int_1^n \frac{1}{u^{\rho+1}} du = \frac{u^{-\rho}}{-\rho} \Big|_{u=1}^n = \frac{1 - n^{-\rho}}{\rho} = \Theta(1)$$

and therefore

$$T(n) = \Theta(n^\rho(1 + \Theta(1))) = \Theta(n^{\lg \phi}).$$

The obvious advantage of the Akra-Bazzi method is that it can solve *almost* any divide-and-conquer recurrence with just a few lines of calculation. (There are a few nasty exceptions like  $T(n) = \sqrt{n} T(\sqrt{n}) + n$  where we have to fall back on recursion trees.) On the other hand, the steps appear to be magic, which makes the method hard to remember, and for *most* divide-and-conquer recurrences, the much simpler recursion tree method is sufficient.

## 5 Linear Recurrences (Annihilators)

Another common class of recurrences, called *linear* recurrences, arises in the context of recursive backtracking algorithms and counting problems. These recurrences express each function value  $f(n)$  as a *linear* combination of a small number of nearby values  $f(n-1), f(n-2), f(n-3), \dots$ . The Fibonacci recurrence is a typical example:

$$F(n) = \begin{cases} 0 & \text{if } n = 0 \\ 1 & \text{if } n = 1 \\ F(n-1) + F(n-2) & \text{otherwise} \end{cases}$$

It turns out that the solution to *any* linear recurrence is a simple combination of polynomial and exponential functions in  $n$ . For example, we can verify by induction that the linear recurrence

$$T(n) = \begin{cases} 1 & \text{if } n = 0 \\ 0 & \text{if } n = 1 \text{ or } n = 2 \\ 3T(n-1) - 8T(n-2) + 4T(n-3) & \text{otherwise} \end{cases}$$

has the closed-form solution  $T(n) = (n-3)2^n + 4$ . First we check the base cases:

$$\begin{aligned} T(0) &= (0-3)2^0 + 4 = 1 & \checkmark \\ T(1) &= (1-3)2^1 + 4 = 0 & \checkmark \\ T(2) &= (2-3)2^2 + 4 = 0 & \checkmark \end{aligned}$$

And now the recursive case:

$$\begin{aligned} T(n) &= 3T(n-1) - 8T(n-2) + 4T(n-3) \\ &= 3((n-4)2^{n-1} + 4) - 8((n-5)2^{n-2} + 4) + 4((n-6)2^{n-3} + 4) \\ &= \left(\frac{3}{2} - \frac{8}{4} + \frac{4}{8}\right)n \cdot 2^n - \left(\frac{12}{2} - \frac{40}{4} + \frac{24}{8}\right)2^n + (2-8+4) \cdot 4 \\ &= (n-3) \cdot 2^n + 4 & \checkmark \end{aligned}$$

But how could we have possibly come up with that solution? In this section, I'll describe a general method for solving linear recurrences that's arguably easier than the induction proof!

### 5.1 Operators

Our technique for solving linear recurrences relies on the theory of *operators*. Operators are higher-order functions, which take one or more functions as input and produce different functions as output. For example, your first two semesters of calculus focus almost exclusively on the *differential* and *integral* operators  $\frac{d}{dx}$  and  $\int dx$ . All the operators we will need are combinations of three elementary building blocks:

- **Sum:**  $(f + g)(n) := f(n) + g(n)$
- **Scale:**  $(\alpha \cdot f)(n) := \alpha \cdot (f(n))$
- **Shift:**  $(Ef)(n) := f(n+1)$

The shift and scale operators are *linear*, which means they can be distributed over sums; for example, for any functions  $f$ ,  $g$ , and  $h$ , we have  $E(f - 3(g - h)) = Ef + (-3)Eg + 3 Eh$ .

We can combine these building blocks to obtain more complex *compound* operators. For example, the compound operator  $E - 2$  is defined by setting  $(E - 2)f := Ef + (-2)f$  for any function  $f$ . We can also apply the shift operator twice:  $(E(Ef))(n) = f(n+2)$ ; we write usually  $E^2f$  as a synonym for  $E(Ef)$ . More generally, for any positive integer  $k$ , the operator  $E^k$  shifts its argument  $k$  times:  $E^k f(n) = f(n+k)$ . Similarly,  $(E - 2)^2$  is shorthand for the operator  $(E - 2)(E - 2)$ , which applies  $(E - 2)$  twice.

For example, here are the results of applying different operators to the function  $f(n) = 2^n$ :

$$\begin{aligned} 2f(n) &= 2 \cdot 2^n = 2^{n+1} \\ 3f(n) &= 3 \cdot 2^n \\ Ef(n) &= 2^{n+1} \\ E^2f(n) &= 2^{n+2} \\ (E - 2)f(n) &= Ef(n) - 2f(n) = 2^{n+1} - 2^{n+1} = 0 \\ (E^2 - 1)f(n) &= E^2f(n) - f(n) = 2^{n+2} - 2^n = 3 \cdot 2^n \end{aligned}$$

These compound operators can be manipulated exactly as though they were polynomials over the ‘variable’  $E$ . In particular, we can ‘factor’ compound operators into ‘products’ of simpler operators, and the order of the factors is unimportant. For example, the compound operators  $E^2 - 3E + 2$  and  $(E - 1)(E - 2)$  are equivalent:

$$\text{Let } g(n) := (E - 2)f(n) = f(n+1) - 2f(n).$$

$$\begin{aligned} \text{Then } (E - 1)(E - 2)f(n) &= (E - 1)g(n) \\ &= g(n+1) - g(n) \\ &= (f(n+2) - 2f(n-1)) - (f(n+1) - 2f(n)) \\ &= f(n+2) - 3f(n+1) + 2f(n) \\ &= (E^2 - 3E + 2)f(n). \quad \checkmark \end{aligned}$$

It is an easy exercise to confirm that  $E^2 - 3E + 2$  is also equivalent to the operator  $(E - 2)(E - 1)$ .

The following table summarizes everything we need to remember about operators.

Operator	Definition
addition	$(f + g)(n) := f(n) + g(n)$
subtraction	$(f - g)(n) := f(n) - g(n)$
multiplication	$(\alpha \cdot f)(n) := \alpha \cdot (f(n))$
shift	$Ef(n) := f(n+1)$
$k$ -fold shift	$E^k f(n) := f(n+k)$
composition	$(X + Y)f := Xf + Yf$ $(X - Y)f := Xf - Yf$ $XYf := X(Yf) = Y(Xf)$
distribution	$X(f + g) = Xf + Xg$

## 5.2 Annihilators

An *annihilator* of a function  $f$  is any nontrivial operator that transforms  $f$  into the zero function. (We can trivially annihilate any function by multiplying it by zero, so as a technical matter, we do not consider

the zero operator to be an annihilator.) Every compound operator we consider annihilates a specific class of functions; conversely, every function composed of polynomial and exponential functions has a unique (minimal) annihilator.

We have already seen that the operator  $(E - 2)$  annihilates the function  $2^n$ . It's not hard to see that the operator  $(E - c)$  annihilates the function  $\alpha \cdot c^n$ , for any constants  $c$  and  $\alpha$ . More generally, the operator  $(E - c)$  annihilates the function  $a^n$  if and only if  $c = a$ :

$$(E - c)a^n = Ea^n - c \cdot a^n = a^{n+1} - c \cdot a^n = (a - c)a^n.$$

Thus,  $(E - 2)$  is essentially the *only* annihilator of the function  $2^n$ .

What about the function  $2^n + 3^n$ ? The operator  $(E - 2)$  annihilates the function  $2^n$ , but leaves the function  $3^n$  unchanged. Similarly,  $(E - 3)$  annihilates  $3^n$  while *negating* the function  $2^n$ . But if we apply *both* operators, we annihilate both terms:

$$\begin{aligned} (E - 2)(2^n + 3^n) &= E(2^n + 3^n) - 2(2^n + 3^n) \\ &= (2^{n+1} + 3^{n+1}) - (2^{n+1} + 2 \cdot 3^n) = 3^n \\ \implies (E - 3)(E - 2)(2^n + 3^n) &= (E - 3)3^n = 0 \end{aligned}$$

In general, for any integers  $a \neq b$ , the operator  $(E - a)(E - b) = (E - b)(E - a) = (E^2 - (a + b)E + ab)$  annihilates any function of the form  $\alpha a^n + \beta b^n$ , but nothing else.

What about the operator  $(E - a)(E - a) = (E - a)^2$ ? It turns out that this operator annihilates all functions of the form  $(\alpha n + \beta)a^n$ :

$$\begin{aligned} (E - a)((\alpha n + \beta)a^n) &= (\alpha(n+1) + \beta)a^{n+1} - a(\alpha n + \beta)a^n \\ &= \alpha a^{n+1} \\ \implies (E - a)^2((\alpha n + \beta)a^n) &= (E - a)(\alpha a^{n+1}) = 0 \end{aligned}$$

More generally, the operator  $(E - a)^d$  annihilates all functions of the form  $p(n) \cdot a^n$ , where  $p(n)$  is a polynomial of degree at most  $d - 1$ . For example,  $(E - 1)^3$  annihilates any polynomial of degree at most 2.

The following table summarizes everything we need to remember about annihilators.

Operator	Functions annihilated
$E - 1$	$\alpha$
$E - a$	$\alpha a^n$
$(E - a)(E - b)$	$\alpha a^n + \beta b^n$ [if $a \neq b$ ]
$(E - a_0)(E - a_1) \cdots (E - a_k)$	$\sum_{i=0}^k \alpha_i a_i^n$ [if $a_i$ distinct]
$(E - 1)^2$	$\alpha n + \beta$
$(E - a)^2$	$(\alpha n + \beta)a^n$
$(E - a)^2(E - b)$	$(\alpha n + \beta)a^b + \gamma b^n$ [if $a \neq b$ ]
$(E - a)^d$	$(\sum_{i=0}^{d-1} \alpha_i n^i)a^n$
If $X$ annihilates $f$ , then $X$ also annihilates $Ef$ .	
If $X$ annihilates both $f$ and $g$ , then $X$ also annihilates $f \pm g$ .	
If $X$ annihilates $f$ , then $X$ also annihilates $\alpha f$ , for any constant $\alpha$ .	
If $X$ annihilates $f$ and $Y$ annihilates $g$ , then $XY$ annihilates $f \pm g$ .	

### 5.3 Annihilating Recurrences

Given a linear recurrence for a function, it's easy to extract an annihilator for that function. For many recurrences, we only need to rewrite the recurrence in operator notation. Once we have an annihilator, we can factor it into operators of the form  $(E - c)$ ; the table on the previous page then gives us a generic solution with some unknown coefficients. If we are given explicit base cases, we can determine the coefficients by examining a few small cases; in general, this involves solving a small system of linear equations. If the base cases are not specified, the generic solution almost always gives us an asymptotic solution. Here is the technique step by step:

1. Write the recurrence in operator form
2. Extract an annihilator for the recurrence
3. Factor the annihilator (if necessary)
4. Extract the *generic solution* from the annihilator
5. Solve for coefficients using base cases (if known)

Here are several examples of the technique in action:

- $r(n) = 5r(n - 1)$ , where  $r(0) = 3$ .

1. We can write the recurrence in operator form as follows:

$$r(n) = 5r(n - 1) \implies r(n + 1) - 5r(n) = 0 \implies (E - 5)r(n) = 0.$$

2. We immediately see that  $(E - 5)$  annihilates the function  $r(n)$ .
3. The annihilator  $(E - 5)$  is already factored.
4. Consulting the annihilator table on the previous page, we find the generic solution  $r(n) = \alpha 5^n$  for some constant  $\alpha$ .
5. The base case  $r(0) = 3$  implies that  $\alpha = 3$ .

We conclude that  $r(n) = 3 \cdot 5^n$ . We can easily verify this closed-form solution by induction:

$$\begin{aligned} r(0) &= 3 \cdot 5^0 = 3 \quad \checkmark && [\text{definition}] \\ r(n) &= 5r(n - 1) && [\text{definition}] \\ &= 5 \cdot (3 \cdot 5^{n-1}) && [\text{induction hypothesis}] \\ &= 5^n \cdot 3 \quad \checkmark && [\text{algebra}] \end{aligned}$$

- **Fibonacci numbers:**  $F(n) = F(n - 1) + F(n - 2)$ , where  $F(0) = 0$  and  $F(1) = 1$ .

1. We can rewrite the recurrence as  $(E^2 - E - 1)F(n) = 0$ .
2. The operator  $E^2 - E - 1$  clearly annihilates  $F(n)$ .
3. The quadratic formula implies that the annihilator  $E^2 - E - 1$  factors into  $(E - \phi)(E - \hat{\phi})$ , where  $\phi = (1 + \sqrt{5})/2 \approx 1.618034$  is the golden ratio and  $\hat{\phi} = (1 - \sqrt{5})/2 = 1 - \phi = -1/\phi \approx -0.618034$ .
4. The annihilator implies that  $F(n) = \alpha\phi^n + \hat{\alpha}\hat{\phi}^n$  for some unknown constants  $\alpha$  and  $\hat{\alpha}$ .

5. The base cases give us two equations in two unknowns:

$$\begin{aligned} F(0) &= 0 = \alpha + \hat{\alpha} \\ F(1) &= 1 = \alpha\phi + \hat{\alpha}\hat{\phi} \end{aligned}$$

Solving this system of equations gives us  $\alpha = 1/(2\phi - 1) = 1/\sqrt{5}$  and  $\hat{\alpha} = -1/\sqrt{5}$ .

We conclude with the following exact closed form for the  $n$ th Fibonacci number:

$$F(n) = \frac{\phi^n - \hat{\phi}^n}{\sqrt{5}} = \frac{1}{\sqrt{5}} \left( \frac{1 + \sqrt{5}}{2} \right)^n - \frac{1}{\sqrt{5}} \left( \frac{1 - \sqrt{5}}{2} \right)^n$$

With all the square roots in this formula, it's quite amazing that Fibonacci numbers are integers. However, if we do all the math correctly, all the square roots cancel out when  $i$  is an integer. (In fact, this is pretty easy to prove using the binomial theorem.)

- **Towers of Hanoi:**  $T(n) = 2T(n - 1) + 1$ , where  $T(0) = 0$ . This is our first example of a *non-homogeneous* recurrence, which means the recurrence has one or more non-recursive terms.

1. We can rewrite the recurrence as  $(E - 2)T(n) = 1$ .
2. The operator  $(E - 2)$  doesn't quite annihilate the function; it leaves a *residue* of 1. But we can annihilate the residue by applying the operator  $(E - 1)$ . Thus, the compound operator  $(E - 1)(E - 2)$  annihilates the function.
3. The annihilator is already factored.
4. The annihilator table gives us the generic solution  $T(n) = \alpha 2^n + \beta$  for some unknown constants  $\alpha$  and  $\beta$ .
5. The base cases give us  $T(0) = 0 = \alpha 2^0 + \beta$  and  $T(1) = 1 = \alpha 2^1 + \beta$ . Solving this system of equations, we find that  $\alpha = 1$  and  $\beta = -1$ .

We conclude that  $T(n) = 2^n - 1$ .

For the remaining examples, I won't explicitly enumerate the steps in the solution.

- **Height-balanced trees:**  $H(n) = H(n - 1) + H(n - 2) + 1$ , where  $H(-1) = 0$  and  $H(0) = 1$ . (Yes, we're starting at  $-1$  instead of  $0$ . So what?)

We can rewrite the recurrence as  $(E^2 - E - 1)H = 1$ . The residue 1 is annihilated by  $(E - 1)$ , so the compound operator  $(E - 1)(E^2 - E - 1)$  annihilates the recurrence. This operator factors into  $(E - 1)(E - \phi)(E - \hat{\phi})$ , where  $\phi = (1 + \sqrt{5})/2$  and  $\hat{\phi} = (1 - \sqrt{5})/2$ . Thus, we get the generic solution  $H(n) = \alpha \cdot \phi^n + \beta + \gamma \cdot \hat{\phi}^n$ , for some unknown constants  $\alpha, \beta, \gamma$  that satisfy the following system of equations:

$$\begin{aligned} H(-1) &= 0 = \alpha\phi^{-1} + \beta + \gamma\hat{\phi}^{-1} = \alpha/\phi + \beta - \gamma/\hat{\phi} \\ H(0) &= 1 = \alpha\phi^0 + \beta + \gamma\hat{\phi}^0 = \alpha + \beta + \gamma \\ H(1) &= 2 = \alpha\phi^1 + \beta + \gamma\hat{\phi}^1 = \alpha\phi + \beta + \gamma\hat{\phi} \end{aligned}$$

Solving this system (using Cramer's rule or Gaussian elimination), we find that  $\alpha = (\sqrt{5} + 2)/\sqrt{5}$ ,  $\beta = -1$ , and  $\gamma = (\sqrt{5} - 2)/\sqrt{5}$ . We conclude that

$$H(n) = \frac{\sqrt{5} + 2}{\sqrt{5}} \left( \frac{1 + \sqrt{5}}{2} \right)^n - 1 + \frac{\sqrt{5} - 2}{\sqrt{5}} \left( \frac{1 - \sqrt{5}}{2} \right)^n.$$

- $T(n) = 3T(n-1) - 8T(n-2) + 4T(n-3)$ , where  $T(0) = 1$ ,  $T(1) = 0$ , and  $T(2) = 0$ . This was our original example of a linear recurrence.

We can rewrite the recurrence as  $(E^3 - 3E^2 + 8E - 4)T = 0$ , so we immediately have an annihilator  $E^3 - 3E^2 + 8E - 4$ . Using high-school algebra, we can factor the annihilator into  $(E-2)^2(E-1)$ , which implies the generic solution  $T(n) = \alpha n2^n + \beta 2^n + \gamma$ . The constants  $\alpha$ ,  $\beta$ , and  $\gamma$  are determined by the base cases:

$$\begin{aligned} T(0) &= 1 = \alpha \cdot 0 \cdot 2^0 + \beta 2^0 + \gamma = \beta + \gamma \\ T(1) &= 0 = \alpha \cdot 1 \cdot 2^1 + \beta 2^1 + \gamma = 2\alpha + 2\beta + \gamma \\ T(2) &= 0 = \alpha \cdot 2 \cdot 2^2 + \beta 2^2 + \gamma = 8\alpha + 4\beta + \gamma \end{aligned}$$

Solving this system of equations, we find that  $\alpha = 1$ ,  $\beta = -3$ , and  $\gamma = 4$ , so  $T(n) = (n-3)2^n + 4$ .

- $T(n) = T(n-1) + 2T(n-2) + 2^n - n^2$

We can rewrite the recurrence as  $(E^2 - E - 2)T(n) = E^2(2^n - n^2)$ . Notice that we had to shift up the non-recursive parts of the recurrence when we expressed it in this form. The operator  $(E-2)(E-1)^3$  annihilates the residue  $2^n - n^2$ , and therefore also annihilates the shifted residue  $E^2(2^n + n^2)$ . Thus, the operator  $(E-2)(E-1)^3(E^2 - E - 2)$  annihilates the entire recurrence. We can factor the quadratic factor into  $(E-2)(E+1)$ , so the annihilator factors into  $(E-2)^2(E-1)^3(E+1)$ . So the generic solution is  $T(n) = \alpha n2^n + \beta 2^n + \gamma n^2 + \delta n + \epsilon + \eta(-1)^n$ . The coefficients  $\alpha$ ,  $\beta$ ,  $\gamma$ ,  $\delta$ ,  $\epsilon$ ,  $\eta$  satisfy a system of six equations determined by the first six function values  $T(0)$  through  $T(5)$ . For almost<sup>2</sup> every set of base cases, we have  $\alpha \neq 0$ , which implies that  $T(n) = \Theta(n2^n)$ .

For a more detailed explanation of the annihilator method, see George Lueker, Some techniques for solving recurrences, *ACM Computing Surveys* 12(4):419-436, 1980.

## 6 Transformations

Sometimes we encounter recurrences that don't fit the structures required for recursion trees or annihilators. In many of those cases, we can transform the recurrence into a more familiar form, by defining a new function in terms of the one we want to solve. There are many different kinds of transformations, but these three are probably the most useful:

- Domain transformation: Define a new function  $S(n) = T(f(n))$  with a simpler recurrence, for some simple function  $f$ .
- Range transformation: Define a new function  $S(n) = f(T(n))$  with a simpler recurrence, for some simple function  $f$ .
- Difference transformation: Simplify the recurrence for  $T(n)$  by considering the difference  $T(n) - T(n-1)$ .

Here are some examples of these transformations in action.

- **Unsimplified Mergesort:**  $T(n) = T(\lceil n/2 \rceil) + T(\lfloor n/2 \rfloor) + \Theta(n)$

When  $n$  is a power of 2, we can simplify the mergesort recurrence to  $T(n) = 2T(n/2) + \Theta(n)$ , which has the solution  $T(n) = \Theta(n \log n)$ . Unfortunately, for other values of  $n$ , this simplified

---

<sup>2</sup>In fact, the only possible solutions with  $\alpha = 0$  have the form  $-2^{n-1} - n^2/2 - 5n/2 + \eta(-1)^n$  for some constant  $\eta$ .

recurrence is incorrect. When  $n$  is odd, then the recurrence calls for us to sort a fractional number of elements! Worse yet, if  $n$  is not a power of 2, we will *never* reach the base case  $T(1) = 1$ .

So we really need to solve the original recurrence. We have no hope of getting an *exact* solution, even if we ignore the  $\Theta()$  in the recurrence; the floors and ceilings will eventually kill us. But we can derive a tight asymptotic solution using a domain transformation—we can rewrite the function  $T(n)$  as a nested function  $S(f(n))$ , where  $f(n)$  is a simple function and the function  $S()$  has a simpler recurrence.

First let's overestimate the time bound, once by pretending that the two subproblem sizes are equal, and again to eliminate the ceiling:

$$T(n) \leq 2T(\lceil n/2 \rceil) + n \leq 2T(n/2 + 1) + n.$$

Now we define a new function  $S(n) = T(n + \alpha)$ , where  $\alpha$  is an unknown constant, chosen so that  $S(n)$  satisfies the Master-Theorem-ready recurrence  $S(n) \leq 2S(n/2) + O(n)$ . To figure out the correct value of  $\alpha$ , we compare two versions of the recurrence for the function  $T(n + \alpha)$ :

$$\begin{aligned} S(n) &\leq 2S(n/2) + O(n) &\implies T(n + \alpha) &\leq 2T(n/2 + \alpha) + O(n) \\ T(n) &\leq 2T(n/2 + 1) + n &\implies T(n + \alpha) &\leq 2T((n + \alpha)/2 + 1) + n + \alpha \end{aligned}$$

For these two recurrences to be equal, we need  $n/2 + \alpha = (n + \alpha)/2 + 1$ , which implies that  $\alpha = 2$ . The Master Theorem now tells us that  $S(n) = O(n \log n)$ , so

$$T(n) = S(n - 2) = O((n - 2) \log(n - 2)) = O(n \log n).$$

A similar argument implies the matching lower bound  $T(n) = \Omega(n \log n)$ . So  $T(n) = \Theta(n \log n)$  after all, just as though we had ignored the floors and ceilings from the beginning!

Domain transformations are useful for removing floors, ceilings, and lower order terms from the arguments of any recurrence that otherwise looks like it ought to fit either the Master Theorem or the recursion tree method. But now that we know this, we don't need to bother grinding through the actual gory details!

- **Ham-Sandwich Trees:**  $T(n) = T(n/2) + T(n/4) + 1$

As we saw earlier, the recursion tree method only gives us the uselessly loose bounds  $\sqrt{n} \ll T(n) \ll n$  for this recurrence, and the recurrence is in the wrong form for annihilators. The authors who discovered ham-sandwich trees (yes, this is a real data structure) solved this recurrence by guessing the solution and giving a complicated induction proof.

But a simple transformation allows us to solve the recurrence in just a few lines. We define a new function  $t(k) = T(2^k)$ , which satisfies the simpler linear recurrence  $t(k) = t(k - 1) + t(k - 2) + 1$ . This recurrence should immediately remind you of Fibonacci numbers. Sure enough, the annihilator method implies the solution  $t(k) = \Theta(\phi^k)$ , where  $\phi = (1 + \sqrt{5})/2$  is the golden ratio. We conclude that

$$T(n) = t(\lg n) = \Theta(\phi^{\lg n}) = \Theta(n^{\lg \phi}) \approx \Theta(n^{0.69424}).$$

Recall that we obtained this same solution earlier using the Akra-Bazzi theorem.

Many other divide-and-conquer recurrences can be similarly transformed into linear recurrences and then solved with annihilators. Consider once more the simplified mergesort recurrence

$T(n) = 2T(n/2) + n$ . The function  $t(k) = T(2^k)$  satisfies the recurrence  $t(k) = 2t(k-1) + 2^k$ . The annihilator method gives us the generic solution  $t(k) = \Theta(k \cdot 2^k)$ , which implies that  $T(n) = t(\lg n) = \Theta(n \log n)$ , just as we expected.

On the other hand, for some recurrences like  $T(n) = T(n/3) + T(2n/3) + n$ , the recursion tree method gives an easy solution, but there's no way to transform the recurrence into a form where we can apply the annihilator method directly.<sup>3</sup>

- **Random Binary Search Trees:**  $T(n) = \frac{1}{4}T(n/4) + \frac{3}{4}T(3n/4) + 1$

This looks like a divide-and-conquer recurrence, so we might be tempted to apply recursion trees, but what does it mean to have a quarter of a child? If we're not comfortable with weighted recursion trees or the Akra-Bazzi theorem, we can instead consider a new function  $U(n) = n \cdot T(n)$ , which satisfies the more palatable recurrence  $U(n) = U(n/4) + U(3n/4) + n$ . As we've already seen, recursion trees imply that  $U(n) = \Theta(n \log n)$ , which immediately implies that  $T(n) = \Theta(\log n)$ .

- **Randomized Quicksort:**  $T(n) = \frac{2}{n} \sum_{k=0}^{n-1} T(k) + n$

This is our first example of a *full history* recurrence; each function value  $T(n)$  is defined in terms of *all* previous function values  $T(k)$  with  $k < n$ . Before we can apply any of our existing techniques, we need to convert this recurrence into an equivalent *limited history* form by shifting and subtracting away common terms. To make this step slightly easier, we first multiply both sides of the recurrence by  $n$  to get rid of the fractions.

$$\begin{aligned} n \cdot T(n) &= 2 \sum_{k=0}^{n-1} T(k) + n^2 && [\text{multiply both sides by } n] \\ (n-1) \cdot T(n-1) &= 2 \sum_{k=0}^{n-2} T(k) + (n-1)^2 && [\text{shift}] \\ nT(n) - (n-1)T(n-1) &= 2T(n-1) + 2n - 1 && [\text{subtract}] \\ T(n) &= \frac{n+1}{n} T(n-1) + 2 - \frac{1}{n} && [\text{simplify}] \end{aligned}$$

We can solve this limited-history recurrence using another functional transformation. We define a new function  $t(n) = T(n)/(n+1)$ , which satisfies the simpler recurrence

$$t(n) = t(n-1) + \frac{2}{n+1} - \frac{1}{n(n+1)},$$

which we can easily unroll into a summation. If we only want an asymptotic solution, we can simplify the final recurrence to  $t(n) = t(n-1) + \Theta(1/n)$ , which unrolls into a very familiar summation:

$$t(n) = \sum_{i=1}^n \Theta(1/i) = \Theta(H_n) = \Theta(\log n).$$

---

<sup>3</sup>However, we can still get a solution via functional transformations as follows. The function  $t(k) = T((3/2)^k)$  satisfies the recurrence  $t(n) = t(n-1) + t(n-\lambda) + (3/2)^k$ , where  $\lambda = \log_{3/2} 3 = 2.709511\dots$ . The **characteristic function** for this recurrence is  $(r^\lambda - r^{\lambda-1} - 1)(r - 3/2)$ , which has a double root at  $r = 3/2$  and nowhere else. Thus,  $t(k) = \Theta(k(3/2)^k)$ , which implies that  $T(n) = t(\log_{3/2} n) = \Theta(n \log n)$ . This line of reasoning is the core of the Akra-Bazzi method.

Finally, substituting  $T(n) = (n + 1)t(n)$  gives us a solution to the original recurrence:  $T(n) = \Theta(n \log n)$ .

## Exercises

- For each of the following recurrences, first *guess* an exact closed-form solution, and then prove your guess is correct. You are free to use any method you want to make your guess—unrolling the recurrence, writing out the first several values, induction proof template, recursion trees, annihilators, transformations, ‘It looks like that other one’, whatever—but please describe your method. All functions are from the non-negative integers to the reals. If it simplifies your solutions, express them in terms of Fibonacci numbers  $F_n$ , harmonic numbers  $H_n$ , binomial coefficients  $\binom{n}{k}$ , factorials  $n!$ , and/or the floor and ceiling functions  $\lfloor x \rfloor$  and  $\lceil x \rceil$ .

(a)  $A(n) = A(n - 1) + 1$ , where  $A(0) = 0$ .

(b)  $B(n) = \begin{cases} 0 & \text{if } n < 5 \\ B(n - 5) + 2 & \text{otherwise} \end{cases}$

(c)  $C(n) = C(n - 1) + 2n - 1$ , where  $C(0) = 0$ .

(d)  $D(n) = D(n - 1) + \binom{n}{2}$ , where  $D(0) = 0$ .

(e)  $E(n) = E(n - 1) + 2^n$ , where  $E(0) = 0$ .

(f)  $F(n) = 3 \cdot F(n - 1)$ , where  $F(0) = 1$ .

(g)  $G(n) = \frac{G(n-1)}{G(n-2)}$ , where  $G(0) = 1$  and  $G(1) = 2$ . [Hint: This is easier than it looks.]

(h)  $H(n) = H(n - 1) + 1/n$ , where  $H(0) = 0$ .

(i)  $I(n) = I(n - 2) + 3/n$ , where  $I(0) = I(1) = 0$ . [Hint: Consider even and odd  $n$  separately.]

(j)  $J(n) = J(n - 1)^2$ , where  $J(0) = 2$ .

(k)  $K(n) = K(\lfloor n/2 \rfloor) + 1$ , where  $K(0) = 0$ .

(l)  $L(n) = L(n - 1) + L(n - 2)$ , where  $L(0) = 2$  and  $L(1) = 1$ .

[Hint: Write the solution in terms of Fibonacci numbers.]

(m)  $M(n) = M(n - 1) \cdot M(n - 2)$ , where  $M(0) = 2$  and  $M(1) = 1$ .

[Hint: Write the solution in terms of Fibonacci numbers.]

(n)  $N(n) = 1 + \sum_{k=1}^n (N(k - 1) + N(n - k))$ , where  $N(0) = 1$ .

(p)  $P(n) = \sum_{k=0}^{n-1} (k \cdot P(k - 1))$ , where  $P(0) = 1$ .

(q)  $Q(n) = \frac{1}{2 - Q(n-1)}$ , where  $Q(0) = 0$ .

(r)  $R(n) = \max_{1 \leq k \leq n} \{R(k - 1) + R(n - k) + n\}$

(s)  $S(n) = \max_{1 \leq k \leq n} \{S(k - 1) + S(n - k) + 1\}$

(t)  $T(n) = \min_{1 \leq k \leq n} \{T(k - 1) + T(n - k) + n\}$

- (u)  $U(n) = \min_{1 \leq k \leq n} \{U(k-1) + U(n-k) + 1\}$
- (v)  $V(n) = \max_{n/3 \leq k \leq 2n/3} \{V(k-1) + V(n-k) + n\}$

2. Use recursion trees or the Akra-Bazzi theorem to solve each of the following recurrences.

- (a)  $A(n) = 2A(n/4) + \sqrt{n}$
- (b)  $B(n) = 2B(n/4) + n$
- (c)  $C(n) = 2C(n/4) + n^2$
- (d)  $D(n) = 3D(n/3) + \sqrt{n}$
- (e)  $E(n) = 3E(n/3) + n$
- (f)  $F(n) = 3F(n/3) + n^2$
- (g)  $G(n) = 4G(n/2) + \sqrt{n}$
- (h)  $H(n) = 4H(n/2) + n$
- (i)  $I(n) = 4I(n/2) + n^2$
- (j)  $J(n) = J(n/2) + J(n/3) + J(n/6) + n$
- (k)  $K(n) = K(n/2) + K(n/3) + K(n/6) + n^2$
- (l)  $L(n) = L(n/15) + L(n/10) + 2L(n/6) + \sqrt{n}$
- \*(m)  $M(n) = 2M(n/3) + 2M(2n/3) + n$
- (n)  $N(n) = \sqrt{2n}N(\sqrt{2n}) + \sqrt{n}$
- (p)  $P(n) = \sqrt{2n}P(\sqrt{2n}) + n$
- (q)  $Q(n) = \sqrt{2n}Q(\sqrt{2n}) + n^2$
- (r)  $R(n) = R(n-3) + 8^n$  — Don't use annihilators!
- (s)  $S(n) = 2S(n-2) + 4^n$  — Don't use annihilators!
- (t)  $T(n) = 4T(n-1) + 2^n$  — Don't use annihilators!

3. Make up a bunch of linear recurrences and then solve them using annihilators.

4. Solve the following recurrences, using any tricks at your disposal.

- (a)  $T(n) = \sum_{i=1}^{\lg n} T(n/2^i) + n$       [Hint: Assume  $n$  is a power of 2.]
- (b) More to come...