

Chapter 13

Randomized Algorithms

The idea that a process can be “random” is not a modern one; we can trace the notion far back into the history of human thought and certainly see its reflections in gambling and the insurance business, each of which reach into ancient times. Yet, while similarly intuitive subjects like geometry and logic have been treated mathematically for several thousand years, the mathematical study of probability is surprisingly young; the first known attempts to seriously formalize it came about in the 1600s. Of course, the history of computer science plays out on a much shorter time scale, and the idea of randomization has been with it since its early days.

Randomization and probabilistic analysis are themes that cut across many areas of computer science, including algorithm design, and when one thinks about random processes in the context of computation, it is usually in one of two distinct ways. One view is to consider the world as behaving randomly: One can consider traditional algorithms that confront randomly generated input. This approach is often termed *average-case analysis*, since we are studying the behavior of an algorithm on an “average” input (subject to some underlying random process), rather than a worst-case input.

A second view is to consider algorithms that behave randomly: The world provides the same worst-case input as always, but we allow our algorithm to make random decisions as it processes the input. Thus the role of randomization in this approach is purely internal to the algorithm and does not require new assumptions about the nature of the input. It is this notion of a *randomized algorithm* that we will be considering in this chapter.

Why might it be useful to design an algorithm that is allowed to make random decisions? A first answer would be to observe that by allowing randomization, we've made our underlying model more powerful. Efficient deterministic algorithms that always yield the correct answer are a special case of efficient randomized algorithms that only need to yield the correct answer with high probability; they are also a special case of randomized algorithms that are always correct, and run efficiently *in expectation*. Even in a worst-case world, an algorithm that does its own "internal" randomization may be able to offset certain worst-case phenomena. So problems that may not have been solvable by efficient deterministic algorithms may still be amenable to randomized algorithms.

But this is not the whole story, and in fact we'll be looking at randomized algorithms for a number of problems where there exist comparably efficient deterministic algorithms. Even in such situations, a randomized approach often exhibits considerable power for further reasons: It may be conceptually much simpler; or it may allow the algorithm to function while maintaining very little internal state or memory of the past. The advantages of randomization seem to increase further as one considers larger computer systems and networks, with many loosely interacting processes—in other words, a *distributed system*. Here random behavior on the part of individual processes can reduce the amount of explicit communication or synchronization that is required; it is often valuable as a tool for *symmetry-breaking* among processes, reducing the danger of contention and "hot spots." A number of our examples will come from settings like this: regulating access to a shared resource, balancing load on multiple processors, or routing packets through a network. Even a small level of comfort with randomized heuristics can give one considerable leverage in thinking about large systems.

A natural worry in approaching the topic of randomized algorithms is that it requires an extensive knowledge of probability. Of course, it's always better to know more rather than less, and some algorithms are indeed based on complex probabilistic ideas. But one further goal of this chapter is to illustrate *how little* underlying probability is really needed in order to understand many of the well-known algorithms in this area. We will see that there is a small set of useful probabilistic tools that recur frequently, and this chapter will try to develop the tools alongside the algorithms. Ultimately, facility with these tools is as valuable as an understanding of the specific algorithms themselves.

13.1 A First Application: Contention Resolution

We begin with a first application of randomized algorithms—contention resolution in a distributed system—that illustrates the general style of analysis

we will be using for many of the algorithms that follow. In particular, it is a chance to work through some basic manipulations involving *events* and their probabilities, analyzing intersections of events using *independence* as well as unions of events using a simple *Union Bound*. For the sake of completeness, we give a brief summary of these concepts in the final section of this chapter (Section 13.15).

The Problem

Suppose we have n processes P_1, P_2, \dots, P_n , each competing for access to a single shared database. We imagine time as being divided into discrete *rounds*. The database has the property that it can be accessed by at most one process in a single round; if two or more processes attempt to access it simultaneously, then all processes are “locked out” for the duration of that round. So, while each process wants to access the database as often as possible, it’s pointless for all of them to try accessing it in every round; then everyone will be perpetually locked out. What’s needed is a way to divide up the rounds among the processes in an equitable fashion, so that all processes get through to the database on a regular basis.

If it is easy for the processes to communicate with one another, then one can imagine all sorts of direct means for resolving the contention. But suppose that the processes can’t communicate with one another at all; how then can they work out a protocol under which they manage to “take turns” in accessing the database?

Designing a Randomized Algorithm

Randomization provides a natural protocol for this problem, which we can specify simply as follows. For some number $p > 0$ that we’ll determine shortly, each process will attempt to access the database in each round with probability p , independently of the decisions of the other processes. So, if exactly one process decides to make the attempt in a given round, it will succeed; if two or more try, then they will all be locked out; and if none try, then the round is in a sense “wasted.” This type of strategy, in which each of a set of identical processes randomizes its behavior, is the core of the *symmetry-breaking* paradigm that we mentioned initially: If all the processes operated in lockstep, repeatedly trying to access the database at the same time, there’d be no progress; but by randomizing, they “smooth out” the contention.

Analyzing the Algorithm

As with many applications of randomization, the algorithm in this case is extremely simple to state; the interesting issue is to analyze its performance.

Defining Some Basic Events When confronted with a probabilistic system like this, a good first step is to write down some basic events and think about their probabilities. Here's a first event to consider. For a given process P_i and a given round t , let $\mathcal{A}[i, t]$ denote the event that P_i attempts to access the database in round t . We know that each process attempts an access in each round with probability p , so the probability of this event, for any i and t , is $\Pr[\mathcal{A}[i, t]] = p$. For every event, there is also a *complementary event*, indicating that the event did not occur; here we have the complementary event $\overline{\mathcal{A}[i, t]}$ that P_i does not attempt to access the database in round t , with probability

$$\Pr[\overline{\mathcal{A}[i, t]}] = 1 - \Pr[\mathcal{A}[i, t]] = 1 - p.$$

Our real concern is whether a process *succeeds* in accessing the database in a given round. Let $\mathcal{S}[i, t]$ denote this event. Clearly, the process P_i must attempt an access in round t in order to succeed. Indeed, succeeding is equivalent to the following: Process P_i attempts to access the database in round t , and each other process *does not* attempt to access the database in round t . Thus $\mathcal{S}[i, t]$ is equal to the intersection of the event $\mathcal{A}[i, t]$ with all the complementary events $\overline{\mathcal{A}[j, t]}$, for $j \neq i$:

$$\mathcal{S}[i, t] = \mathcal{A}[i, t] \cap \left(\bigcap_{j \neq i} \overline{\mathcal{A}[j, t]} \right).$$

All the events in this intersection are independent, by the definition of the contention-resolution protocol. Thus, to get the probability of $\mathcal{S}[i, t]$, we can multiply the probabilities of all the events in the intersection:

$$\Pr[\mathcal{S}[i, t]] = \Pr[\mathcal{A}[i, t]] \cdot \prod_{j \neq i} \Pr[\overline{\mathcal{A}[j, t]}] = p(1 - p)^{n-1}.$$

We now have a nice, closed-form expression for the probability that P_i succeeds in accessing the database in round t ; we can now ask how to set p so that this success probability is maximized. Observe first that the success probability is 0 for the extreme cases $p = 0$ and $p = 1$ (these correspond to the extreme case in which processes never bother attempting, and the opposite extreme case in which every process tries accessing the database in every round, so that everyone is locked out). The function $f(p) = p(1 - p)^{n-1}$ is positive for values of p strictly between 0 and 1, and its derivative $f'(p) = (1 - p)^{n-1} - (n - 1)p(1 - p)^{n-2}$ has a single zero at the value $p = 1/n$, where the maximum is achieved. Thus we can maximize the success probability by setting $p = 1/n$. (Notice that $p = 1/n$ is a natural intuitive choice as well, if one wants exactly one process to attempt an access in any round.)

When we set $p = 1/n$, we get $\Pr[\mathcal{S}[i, t]] = \frac{1}{n} \left(1 - \frac{1}{n}\right)^{n-1}$. It's worth getting a sense for the asymptotic value of this expression, with the help of the following extremely useful fact from basic calculus.

(13.1)

- (a) The function $\left(1 - \frac{1}{n}\right)^n$ converges monotonically from $\frac{1}{4}$ up to $\frac{1}{e}$ as n increases from 2.
- (b) The function $\left(1 - \frac{1}{n}\right)^{n-1}$ converges monotonically from $\frac{1}{2}$ down to $\frac{1}{e}$ as n increases from 2.

Using (13.1), we see that $1/(en) \leq \Pr[\mathcal{S}[i, t]] \leq 1/(2n)$, and hence $\Pr[\mathcal{S}[i, t]]$ is asymptotically equal to $\Theta(1/n)$.

Waiting for a Particular Process to Succeed Let's consider this protocol with the optimal value $p = 1/n$ for the access probability. Suppose we are interested in how long it will take process P_i to succeed in accessing the database at least once. We see from the earlier calculation that the probability of its succeeding in any one round is not very good, if n is reasonably large. How about if we consider multiple rounds?

Let $\mathcal{F}[i, t]$ denote the “failure event” that process P_i does not succeed in *any* of the rounds 1 through t . This is clearly just the intersection of the complementary events $\overline{\mathcal{S}[i, r]}$ for $r = 1, 2, \dots, t$. Moreover, since each of these events is independent, we can compute the probability of $\mathcal{F}[i, t]$ by multiplication:

$$\Pr[\mathcal{F}[i, t]] = \Pr\left[\bigcap_{r=1}^t \overline{\mathcal{S}[i, r]}\right] = \prod_{r=1}^t \Pr[\overline{\mathcal{S}[i, r]}] = \left[1 - \frac{1}{n} \left(1 - \frac{1}{n}\right)^{n-1}\right]^t.$$

This calculation does give us the value of the probability; but at this point, we're in danger of ending up with some extremely complicated-looking expressions, and so it's important to start thinking asymptotically. Recall that the probability of success was $\Theta(1/n)$ after one round; specifically, it was bounded between $1/(en)$ and $1/(2n)$. Using the expression above, we have

$$\Pr[\mathcal{F}[i, t]] = \prod_{r=1}^t \Pr[\overline{\mathcal{S}[i, r]}] \leq \left(1 - \frac{1}{en}\right)^t.$$

Now we notice that if we set $t = en$, then we have an expression that can be plugged directly into (13.1). Of course en will not be an integer; so we can take $t = \lceil en \rceil$ and write

$$\Pr[\mathcal{F}[i, t]] \leq \left(1 - \frac{1}{en}\right)^{\lceil en \rceil} \leq \left(1 - \frac{1}{en}\right)^{en} \leq \frac{1}{e}.$$

This is a very compact and useful asymptotic statement: The probability that process P_i does not succeed in any of rounds 1 through $\lceil en \rceil$ is upper-bounded by the constant e^{-1} , independent of n . Now, if we increase t by some fairly small factors, the probability that P_i does not succeed in any of rounds 1 through t drops precipitously: If we set $t = \lceil en \rceil \cdot (c \ln n)$, then we have

$$\Pr [\mathcal{F}[i, t]] \leq \left(1 - \frac{1}{en}\right)^t = \left(\left(1 - \frac{1}{en}\right)^{\lceil en \rceil}\right)^{c \ln n} \leq e^{-c \ln n} = n^{-c}.$$

So, asymptotically, we can view things as follows. After $\Theta(n)$ rounds, the probability that P_i has not yet succeeded is bounded by a constant; and between then and $\Theta(n \ln n)$, this probability drops to a quantity that is quite small, bounded by an inverse polynomial in n .

Waiting for All Processes to Get Through Finally, we're in a position to ask the question that was implicit in the overall setup: How many rounds must elapse before there's a high probability that all processes will have succeeded in accessing the database at least once?

To address this, we say that the protocol *fails* after t rounds if some process has not yet succeeded in accessing the database. Let \mathcal{F}_t denote the event that the protocol fails after t rounds; the goal is to find a reasonably small value of t for which $\Pr [\mathcal{F}_t]$ is small.

The event \mathcal{F}_t occurs if and only if one of the events $\mathcal{F}[i, t]$ occurs; so we can write

$$\mathcal{F}_t = \bigcup_{i=1}^n \mathcal{F}[i, t].$$

Previously, we considered intersections of independent events, which were very simple to work with; here, by contrast, we have a union of events that are not independent. Probabilities of unions like this can be very hard to compute exactly, and in many settings it is enough to analyze them using a simple *Union Bound*, which says that the probability of a union of events is upper-bounded by the sum of their individual probabilities:

(13.2) (The Union Bound) *Given events $\mathcal{E}_1, \mathcal{E}_2, \dots, \mathcal{E}_n$, we have*

$$\Pr \left[\bigcup_{i=1}^n \mathcal{E}_i \right] \leq \sum_{i=1}^n \Pr [\mathcal{E}_i].$$

Note that this is not an equality; but the upper bound is good enough when, as here, the union on the left-hand side represents a “bad event” that

we're trying to avoid, and we want a bound on its probability in terms of constituent “bad events” on the right-hand side.

For the case at hand, recall that $\mathcal{F}_t = \bigcup_{i=1}^n \mathcal{F}[i, t]$, and so

$$\Pr [\mathcal{F}_t] \leq \sum_{i=1}^n \Pr [\mathcal{F}[i, t]].$$

The expression on the right-hand side is a sum of n terms, each with the same value; so to make the probability of \mathcal{F}_t small, we need to make each of the terms on the right significantly smaller than $1/n$. From our earlier discussion, we see that choosing $t = \Theta(n)$ will not be good enough, since then each term on the right is only bounded by a constant. If we choose $t = \lceil en \rceil \cdot (c \ln n)$, then we have $\Pr [\mathcal{F}[i, t]] \leq n^{-c}$ for each i , which is what we want. Thus, in particular, taking $t = 2\lceil en \rceil \ln n$ gives us

$$\Pr [\mathcal{F}_t] \leq \sum_{i=1}^n \Pr [\mathcal{F}[i, t]] \leq n \cdot n^{-2} = n^{-1},$$

and so we have shown the following.

(13.3) *With probability at least $1 - n^{-1}$, all processes succeed in accessing the database at least once within $t = 2\lceil en \rceil \ln n$ rounds.*

An interesting observation here is that if we had chosen a value of t equal to $qn \ln n$ for a very small value of q (rather than the coefficient $2e$ that we actually used), then we would have gotten an upper bound for $\Pr [\mathcal{F}[i, t]]$ that was larger than n^{-1} , and hence a corresponding upper bound for the overall failure probability $\Pr [\mathcal{F}_t]$ that was larger than 1—in other words, a completely worthless bound. Yet, as we saw, by choosing larger and larger values for the coefficient q , we can drive the upper bound on $\Pr [\mathcal{F}_t]$ down to n^{-c} for any constant c we want; and this is really a very tiny upper bound. So, in a sense, all the “action” in the Union Bound takes place rapidly in the period when $t = \Theta(n \ln n)$; as we vary the hidden constant inside the $\Theta(\cdot)$, the Union Bound goes from providing no information to giving an extremely strong upper bound on the probability.

We can ask whether this is simply an artifact of using the Union Bound for our upper bound, or whether it's intrinsic to the process we're observing. Although we won't do the (somewhat messy) calculations here, one can show that when t is a small constant times $n \ln n$, there really is a sizable probability that some process has not yet succeeded in accessing the database. So a rapid falling-off in the value of $\Pr [\mathcal{F}_t]$ genuinely does happen over the range $t = \Theta(n \ln n)$. For this problem, as in many problems of this flavor, we're

really identifying the asymptotically “correct” value of t despite our use of the seemingly weak Union Bound.

13.2 Finding the Global Minimum Cut

Randomization naturally suggested itself in the previous example, since we were assuming a model with many processes that could not directly communicate. We now look at a problem on graphs for which a randomized approach comes as somewhat more of a surprise, since it is a problem for which perfectly reasonable deterministic algorithms exist as well.



The Problem

Given an undirected graph $G = (V, E)$, we define a *cut* of G to be a partition of V into two non-empty sets A and B . Earlier, when we looked at network flows, we worked with the closely related definition of an *s-t cut*: there, given a directed graph $G = (V, E)$ with distinguished source and sink nodes s and t , an *s-t cut* was defined to be a partition of V into sets A and B such that $s \in A$ and $t \in B$. Our definition now is slightly different, since the underlying graph is now undirected and there is no source or sink.

For a cut (A, B) in an undirected graph G , the *size* of (A, B) is the number of edges with one end in A and the other in B . A *global minimum cut* (or “global min-cut” for short) is a cut of minimum size. The term *global* here is meant to connote that any cut of the graph is allowed; there is no source or sink. Thus the global min-cut is a natural “robustness” parameter; it is the smallest number of edges whose deletion disconnects the graph. We first check that network flow techniques are indeed sufficient to find a global min-cut.

(13.4) *There is a polynomial-time algorithm to find a global min-cut in an undirected graph G .*

Proof. We start from the similarity between cuts in undirected graphs and *s-t* cuts in directed graphs, and with the fact that we know how to find the latter optimally.

So given an undirected graph $G = (V, E)$, we need to transform it so that there are directed edges and there is a source and sink. We first replace every undirected edge $e = (u, v) \in E$ with two oppositely oriented directed edges, $e' = (u, v)$ and $e'' = (v, u)$, each of capacity 1. Let G' denote the resulting directed graph.

Now suppose we pick two arbitrary nodes $s, t \in V$, and find the minimum *s-t* cut in G' . It is easy to check that if (A, B) is this minimum cut in G' , then (A, B) is also a cut of minimum size in G among all those that separate s from t . But we know that the global min-cut in G must separate s from *something*,

since both sides A and B are nonempty, and s belongs to only one of them. So we fix any $s \in V$ and compute the minimum s - t cut in G' for every other node $t \in V - \{s\}$. This is $n - 1$ directed minimum-cut computations, and the best among these will be a global min-cut of G . ■

The algorithm in (13.4) gives the strong impression that finding a global min-cut in an undirected graph is in some sense a *harder* problem than finding a minimum s - t cut in a flow network, as we had to invoke a subroutine for the latter problem $n - 1$ times in our method for solving the former. But it turns out that this is just an illusion. A sequence of increasingly simple algorithms in the late 1980s and early 1990s showed that global min-cuts in undirected graphs could actually be computed just as efficiently as s - t cuts or even more so, and by techniques that didn't require augmenting paths or even a notion of flow. The high point of this line of work came with David Karger's discovery in 1992 of the Contraction Algorithm, a randomized method that is qualitatively simpler than all previous algorithms for global min-cuts. Indeed, it is sufficiently simple that, on a first impression, it is very hard to believe that it actually works.

Designing the Algorithm

Here we describe the Contraction Algorithm in its simplest form. This version, while it runs in polynomial time, is not among the most efficient algorithms for global min-cuts. However, subsequent optimizations to the algorithm have given it a much better running time.

The Contraction Algorithm works with a connected *multigraph* $G = (V, E)$; this is an undirected graph that is allowed to have multiple “parallel” edges between the same pair of nodes. It begins by choosing an edge $e = (u, v)$ of G uniformly at random and *contracting* it, as shown in Figure 13.1. This means we produce a new graph G' in which u and v have been identified into a single new node w ; all other nodes keep their identity. Edges that had one end equal to u and the other equal to v are deleted from G' . Each other edge e is preserved in G' , but if one of its ends was equal to u or v , then this end is updated to be equal to the new node w . Note that, even if G had at most one edge between any two nodes, G' may end up with parallel edges.

The Contraction Algorithm then continues recursively on G' , choosing an edge uniformly at random and contracting it. As these recursive calls proceed, the constituent vertices of G' should be viewed as *supernodes*: Each supernode w corresponds to the subset $S(w) \subseteq V$ that has been “swallowed up” in the contractions that produced w . The algorithm terminates when it reaches a graph G' that has only two supernodes v_1 and v_2 (presumably with a number of parallel edges between them). Each of these supernodes v_i has a corresponding subset $S(v_i) \subseteq V$ consisting of the nodes that have been

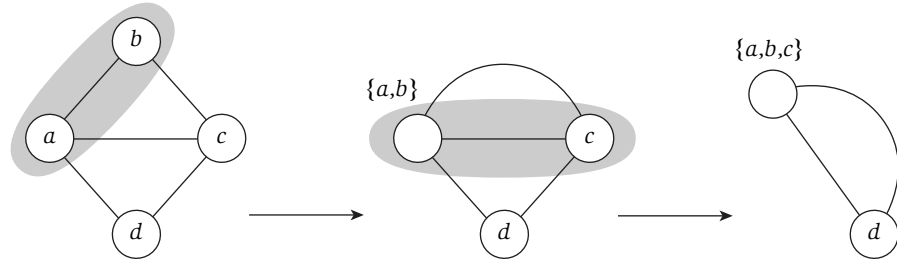


Figure 13.1 The Contraction Algorithm applied to a four-node input graph.

contracted into it, and these two sets $S(v_1)$ and $S(v_2)$ form a partition of V . We output $(S(v_1), S(v_2))$ as the cut found by the algorithm.

```

The Contraction Algorithm applied to a multigraph  $G = (V, E)$ :
  For each node  $v$ , we will record
    the set  $S(v)$  of nodes that have been contracted into  $v$ 
  Initially  $S(v) = \{v\}$  for each  $v$ 
  If  $G$  has two nodes  $v_1$  and  $v_2$ , then return the cut  $(S(v_1), S(v_2))$ 
  Else choose an edge  $e = (u, v)$  of  $G$  uniformly at random
    Let  $G'$  be the graph resulting from the contraction of  $e$ ,
      with a new node  $z_{uv}$  replacing  $u$  and  $v$ 
    Define  $S(z_{uv}) = S(u) \cup S(v)$ 
    Apply the Contraction Algorithm recursively to  $G'$ 
  Endif

```

Analyzing the Algorithm

The algorithm is making random choices, so there is some probability that it will succeed in finding a global min-cut and some probability that it won't. One might imagine at first that the probability of success is exponentially small. After all, there are exponentially many possible cuts of G ; what's favoring the minimum cut in the process? But we'll show first that, in fact, the success probability is only polynomially small. It will then follow that by running the algorithm a polynomial number of times and returning the best cut found in any run, we can actually produce a global min-cut with high probability.

(13.5) *The Contraction Algorithm returns a global min-cut of G with probability at least $1/\binom{n}{2}$.*

Proof. We focus on a global min-cut (A, B) of G and suppose it has size k ; in other words, there is a set F of k edges with one end in A and the other

in B . We want to give a lower bound on the probability that the Contraction Algorithm returns the cut (A, B) .

Consider what could go wrong in the first step of the Contraction Algorithm: The problem would be if an edge in F were contracted. For then, a node of A and a node of B would get thrown together in the same supernode, and (A, B) could not be returned as the output of the algorithm. Conversely, if an edge not in F is contracted, then there is still a chance that (A, B) could be returned.

So what we want is an upper bound on the probability that an edge in F is contracted, and for this we need a lower bound on the size of E . Notice that if any node v had degree less than k , then the cut $(\{v\}, V - \{v\})$ would have size less than k , contradicting our assumption that (A, B) is a global min-cut. Thus every node in G has degree at least k , and so $|E| \geq \frac{1}{2}kn$. Hence the probability that an edge in F is contracted is at most

$$\frac{k}{\frac{1}{2}kn} = \frac{2}{n}.$$

Now consider the situation after j iterations, when there are $n - j$ supernodes in the current graph G' , and suppose that no edge in F has been contracted yet. Every cut of G' is a cut of G , and so there are at least k edges incident to every supernode of G' . Thus G' has at least $\frac{1}{2}k(n - j)$ edges, and so the probability that an edge of F is contracted in the next iteration $j + 1$ is at most

$$\frac{k}{\frac{1}{2}k(n - j)} = \frac{2}{n - j}.$$

The cut (A, B) will actually be returned by the algorithm if no edge of F is contracted in any of iterations $1, 2, \dots, n - 2$. If we write \mathcal{E}_j for the event that an edge of F is not contracted in iteration j , then we have shown $\Pr[\mathcal{E}_1] \geq 1 - 2/n$ and $\Pr[\mathcal{E}_{j+1} | \mathcal{E}_1 \cap \mathcal{E}_2 \cdots \cap \mathcal{E}_j] \geq 1 - 2/(n - j)$. We are interested in lower-bounding the quantity $\Pr[\mathcal{E}_1 \cap \mathcal{E}_2 \cdots \cap \mathcal{E}_{n-2}]$, and we can check by unwinding the formula for conditional probability that this is equal to

$$\begin{aligned} & \Pr[\mathcal{E}_1] \cdot \Pr[\mathcal{E}_2 | \mathcal{E}_1] \cdots \Pr[\mathcal{E}_{j+1} | \mathcal{E}_1 \cap \mathcal{E}_2 \cdots \cap \mathcal{E}_j] \cdots \Pr[\mathcal{E}_{n-2} | \mathcal{E}_1 \cap \mathcal{E}_2 \cdots \cap \mathcal{E}_{n-3}] \\ & \geq \left(1 - \frac{2}{n}\right) \left(1 - \frac{2}{n-1}\right) \cdots \left(1 - \frac{2}{n-j}\right) \cdots \left(1 - \frac{2}{3}\right) \\ & = \left(\frac{n-2}{n}\right) \left(\frac{n-3}{n-1}\right) \left(\frac{n-4}{n-2}\right) \cdots \left(\frac{2}{4}\right) \left(\frac{1}{3}\right) \\ & = \frac{2}{n(n-1)} = \binom{n}{2}^{-1}. \quad \blacksquare \end{aligned}$$

So we now know that a single run of the Contraction Algorithm fails to find a global min-cut with probability at most $(1 - 1/\binom{n}{2})$. This number is very close to 1, of course, but we can amplify our probability of success simply by repeatedly running the algorithm, with independent random choices, and taking the best cut we find. By fact (13.1), if we run the algorithm $\binom{n}{2}$ times, then the probability that we fail to find a global min-cut in any run is at most

$$\left(1 - 1/\binom{n}{2}\right)^{\binom{n}{2}} \leq \frac{1}{e}.$$

And it's easy to drive the failure probability below $1/e$ with further repetitions: If we run the algorithm $\binom{n}{2} \ln n$ times, then the probability we fail to find a global min-cut is at most $e^{-\ln n} = 1/n$.

The overall running time required to get a high probability of success is polynomial in n , since each run of the Contraction Algorithm takes polynomial time, and we run it a polynomial number of times. Its running time will be fairly large compared with the best network flow techniques, since we perform $\Theta(n^2)$ independent runs and each takes at least $\Omega(m)$ time. We have chosen to describe this version of the Contraction Algorithm since it is the simplest and most elegant; it has been shown that some clever optimizations to the way in which multiple runs are performed can improve the running time considerably.

Further Analysis: The Number of Global Minimum Cuts

The analysis of the Contraction Algorithm provides a surprisingly simple answer to the following question: Given an undirected graph $G = (V, E)$ on n nodes, what is the maximum number of global min-cuts it can have (as a function of n)?

For a directed flow network, it's easy to see that the number of minimum s - t cuts can be exponential in n . For example, consider a directed graph with nodes $s, t, v_1, v_2, \dots, v_n$, and unit-capacity edges (s, v_i) and (v_i, t) for each i . Then s together with any subset of $\{v_1, v_2, \dots, v_n\}$ will constitute the source side of a minimum cut, and so there are 2^n minimum s - t cuts.

But for global min-cuts in an undirected graph, the situation looks quite different. If one spends some time trying out examples, one finds that the n -node cycle has $\binom{n}{2}$ global min-cuts (obtained by cutting any two edges), and it is not clear how to construct an undirected graph with more.

We now show how the analysis of the Contraction Algorithm settles this question immediately, establishing that the n -node cycle is indeed an extreme case.

(13.6) An undirected graph $G = (V, E)$ on n nodes has at most $\binom{n}{2}$ global min-cuts.

Proof. The key is that the proof of (13.5) actually established more than was claimed. Let G be a graph, and let C_1, \dots, C_r denote all its global min-cuts. Let \mathcal{E}_i denote the event that C_i is returned by the Contraction Algorithm, and let $\mathcal{E} = \cup_{i=1}^r \mathcal{E}_i$ denote the event that the algorithm returns any global min-cut.

Then, although (13.5) simply asserts that $\Pr[\mathcal{E}] \geq 1/\binom{n}{2}$, its proof actually shows that for each i , we have $\Pr[\mathcal{E}_i] \geq 1/\binom{n}{2}$. Now each pair of events \mathcal{E}_i and \mathcal{E}_j are disjoint—since only one cut is returned by any given run of the algorithm—so by the Union Bound for disjoint events (13.49), we have

$$\Pr[\mathcal{E}] = \Pr\left[\cup_{i=1}^r \mathcal{E}_i\right] = \sum_{i=1}^r \Pr[\mathcal{E}_i] \geq r / \binom{n}{2}.$$

But clearly $\Pr[\mathcal{E}] \leq 1$, and so we must have $r \leq \binom{n}{2}$. ■

13.3 Random Variables and Their Expectations

Thus far our analysis of randomized algorithms and processes has been based on identifying certain “bad events” and bounding their probabilities. This is a qualitative type of analysis, in the sense that the algorithm either succeeds or it doesn’t. A more quantitative style of analysis would consider certain parameters associated with the behavior of the algorithm—for example, its running time, or the quality of the solution it produces—and seek to determine the *expected* size of these parameters over the random choices made by the algorithm. In order to make such analysis possible, we need the fundamental notion of a *random variable*.

Given a probability space, a random variable X is a function from the underlying sample space to the natural numbers, such that for each natural number j , the set $X^{-1}(j)$ of all sample points taking the value j is an event. Thus we can write $\Pr[X = j]$ as loose shorthand for $\Pr[X^{-1}(j)]$; it is because we can ask about X ’s probability of taking a given value that we think of it as a “random variable.”

Given a random variable X , we are often interested in determining its *expectation*—the “average value” assumed by X . We define this as

$$E[X] = \sum_{j=0}^{\infty} j \cdot \Pr[X = j],$$

declaring this to have the value ∞ if the sum diverges. Thus, for example, if X takes each of the values in $\{1, 2, \dots, n\}$ with probability $1/n$, then $E[X] = 1(1/n) + 2(1/n) + \dots + n(1/n) = \binom{n+1}{2}/n = (n+1)/2$.

Example: Waiting for a First Success

Here's a more useful example, in which we see how an appropriate random variable lets us talk about something like the "running time" of a simple random process. Suppose we have a coin that comes up **heads** with probability $p > 0$, and **tails** with probability $1 - p$. Different flips of the coin have independent outcomes. If we flip the coin until we first get a **heads**, what's the expected number of flips we will perform? To answer this, we let X denote the random variable equal to the number of flips performed. For $j > 0$, we have $\Pr[X = j] = (1 - p)^{j-1}p$: in order for the process to take exactly j steps, the first $j - 1$ flips must come up **tails**, and the j^{th} must come up **heads**. Now, applying the definition, we have

$$\begin{aligned} E[X] &= \sum_{j=0}^{\infty} j \cdot \Pr[X = j] = \sum_{j=1}^{\infty} j(1 - p)^{j-1}p = \frac{p}{1 - p} \sum_{j=1}^{\infty} j(1 - p)^j \\ &= \frac{p}{1 - p} \cdot \frac{(1 - p)}{p^2} = \frac{1}{p}. \end{aligned}$$

Thus we get the following intuitively sensible result.

(13.7) *If we repeatedly perform independent trials of an experiment, each of which succeeds with probability $p > 0$, then the expected number of trials we need to perform until the first success is $1/p$.*

Linearity of Expectation

In Sections 13.1 and 13.2, we broke events down into unions of much simpler events, and worked with the probabilities of these simpler events. This is a powerful technique when working with random variables as well, and it is based on the principle of *linearity of expectation*.

(13.8) *Linearity of Expectation. Given two random variables X and Y defined over the same probability space, we can define $X + Y$ to be the random variable equal to $X(\omega) + Y(\omega)$ on a sample point ω . For any X and Y , we have*

$$E[X + Y] = E[X] + E[Y].$$

We omit the proof, which is not difficult. Much of the power of (13.8) comes from the fact that it applies to the sum of *any* random variables; no restrictive assumptions are needed. As a result, if we need to compute the

expectation of a complicated random variable X , we can first write it as a sum of simpler random variables $X = X_1 + X_2 + \cdots + X_n$, compute each $E[X_i]$, and then determine $E[X] = \sum E[X_i]$. We now look at some examples of this principle in action.

Example: Guessing Cards

Memoryless Guessing To amaze your friends, you have them shuffle a deck of 52 cards and then turn over one card at a time. Before each card is turned over, you predict its identity. Unfortunately, you don't have any particular psychic abilities—and you're not so good at remembering what's been turned over already—so your strategy is simply to guess a card uniformly at random from the full deck each time. On how many predictions do you expect to be correct?

Let's work this out for the more general setting in which the deck has n distinct cards, using X to denote the random variable equal to the number of correct predictions. A surprisingly effortless way to compute X is to define the random variable X_i , for $i = 1, 2, \dots, n$, to be equal to 1 if the i^{th} prediction is correct, and 0 otherwise. Notice that $X = X_1 + X_2 + \cdots + X_n$, and

$$E[X_i] = 0 \cdot \Pr[X_i = 0] + 1 \cdot \Pr[X_i = 1] = \Pr[X_i = 1] = \frac{1}{n}.$$

It's worth pausing to note a useful fact that is implicitly demonstrated by the above calculation: If Z is any random variable that only takes the values 0 or 1, then $E[Z] = \Pr[Z = 1]$.

Since $E[X_i] = \frac{1}{n}$ for each i , we have

$$E[X] = \sum_{i=1}^n E[X_i] = n \left(\frac{1}{n} \right) = 1.$$

Thus we have shown the following.

(13.9) *The expected number of correct predictions under the memoryless guessing strategy is 1, independent of n .*

Trying to compute $E[X]$ directly from the definition $\sum_{j=0}^{\infty} j \cdot \Pr[X = j]$ would be much more painful, since it would involve working out a much more elaborate summation. A significant amount of complexity is hidden away in the seemingly innocuous statement of (13.8).

Guessing with Memory Now let's consider a second scenario. Your psychic abilities have not developed any further since last time, but you have become very good at remembering which cards have already been turned over. Thus, when you predict the next card now, you only guess uniformly from among

the cards *not yet seen*. How many correct predictions do you expect to make with this strategy?

Again, let the random variable X_i take the value 1 if the i^{th} prediction is correct, and 0 otherwise. In order for the i^{th} prediction to be correct, you need only guess the correct one out of $n - i + 1$ remaining cards; hence

$$E[X_i] = \Pr[X_i = 1] = \frac{1}{n - i + 1},$$

and so we have

$$\Pr[X] = \sum_{i=1}^n E[X_i] = \sum_{i=1}^n \frac{1}{n - i + 1} = \sum_{i=1}^n \frac{1}{i}.$$

This last expression $\sum_{i=1}^n \frac{1}{i} = 1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n}$ is the *harmonic number* $H(n)$, and it is something that has come up in each of the previous two chapters. In particular, we showed in Chapter 11 that $H(n)$, as a function of n , closely shadows the value $\int_1^{n+1} \frac{1}{x} dx = \ln(n+1)$. For our purposes here, we restate the basic bound on $H(n)$ as follows.

(13.10) $\ln(n+1) < H(n) < 1 + \ln n$, and more loosely, $H(n) = \Theta(\log n)$.

Thus, once you are able to remember the cards you've already seen, the expected number of correct predictions increases significantly above 1.

(13.11) *The expected number of correct predictions under the guessing strategy with memory is $H(n) = \Theta(\log n)$.*

Example: Collecting Coupons

Before moving on to more sophisticated applications, let's consider one more basic example in which linearity of expectation provides significant leverage.

Suppose that a certain brand of cereal includes a free coupon in each box. There are n different types of coupons. As a regular consumer of this brand, how many boxes do you expect to buy before finally getting a coupon of each type?

Clearly, at least n boxes are needed; but it would be sort of surprising if you actually had all n types of coupons by the time you'd bought n boxes. As you collect more and more different types, it will get less and less likely that a new box has a type of coupon you haven't seen before. Once you have $n - 1$ of the n different types, there's only a probability of $1/n$ that a new box has the missing type you need.

Here's a way to work out the expected time exactly. Let X be the random variable equal to the number of boxes you buy until you first have a coupon

of each type. As in our previous examples, this is a reasonably complicated random variable to think about, and we'd like to write it as a sum of simpler random variables. To think about this, let's consider the following natural idea: The coupon-collecting process *makes progress* whenever you buy a box of cereal containing a type of coupon you haven't seen before. Thus the goal of the process is really to make progress n times. Now, at a given point in time, what is the probability that you make progress in the next step? This depends on how many different types of coupons you already have. If you have j types, then the probability of making progress in the next step is $(n - j)/n$: Of the n types of coupons, $n - j$ allow you to make progress. Since the probability varies depending on the number of different types of coupons we have, this suggests a natural way to break down X into simpler random variables, as follows.

Let's say that the coupon-collecting process is in *phase* j when you've already collected j different types of coupons and are waiting to get a new type. When you see a new type of coupon, phase j ends and phase $j + 1$ begins. Thus we start in phase 0, and the whole process is done at the end of phase $n - 1$. Let X_j be the random variable equal to the number of steps you spend in phase j . Then $X = X_0 + X_1 + \cdots + X_{n-1}$, and so it is enough to work out $E[X_j]$ for each j .

$$(13.12) \quad E[X_j] = n/(n - j).$$

Proof. In each step of phase j , the phase ends immediately if and only if the coupon you get next is one of the $n - j$ types you haven't seen before. Thus, in phase j , you are really just waiting for an event of probability $(n - j)/n$ to occur, and so, by (13.7), the expected length of phase j is $E[X_j] = n/(n - j)$. ■

Using this, linearity of expectation gives us the overall expected time.

(13.13) *The expected time before all n types of coupons are collected is $E[X] = nH(n) = \Theta(n \log n)$.*

Proof. By linearity of expectation, we have

$$E[X] = \sum_{j=0}^{n-1} E[X_j] = \sum_{j=0}^{n-1} \frac{n}{n-j} = n \sum_{j=0}^{n-1} \frac{1}{n-j} = n \sum_{i=1}^n \frac{1}{i} = nH(n).$$

By (13.10), we know this is asymptotically equal to $\Theta(n \log n)$. ■

It is interesting to compare the dynamics of this process to one's intuitive view of it. Once $n - 1$ of the n types of coupons are collected, you expect to

buy n more boxes of cereal before you see the final type. In the meantime, you keep getting coupons you've already seen before, and you might conclude that this final type is "the rare one." But in fact it's just as likely as all the others; it's simply that the final one, whichever it turns out to be, is likely to take a long time to get.

A Final Definition: Conditional Expectation

We now discuss one final, very useful notion concerning random variables that will come up in some of the subsequent analyses. Just as one can define the conditional probability of one event given another, one can analogously define the expectation of a random variable conditioned on a certain event. Suppose we have a random variable X and an event \mathcal{E} of positive probability. Then we define the *conditional expectation* of X , given \mathcal{E} , to be the expected value of X computed only over the part of the sample space corresponding to \mathcal{E} . We denote this quantity by $E[X | \mathcal{E}]$. This simply involves replacing the probabilities $\Pr[X = j]$ in the definition of the expectation with conditional probabilities:

$$E[X | \mathcal{E}] = \sum_{j=0}^{\infty} j \cdot \Pr[X = j | \mathcal{E}].$$

13.4 A Randomized Approximation Algorithm for MAX 3-SAT

In the previous section, we saw a number of ways in which linearity of expectation can be used to analyze a randomized process. We now describe an application of this idea to the design of an approximation algorithm. The problem we consider is a variation of the 3-SAT Problem, and we will see that one consequence of our randomized approximation algorithm is a surprisingly strong general statement about 3-SAT that on its surface seems to have nothing to do with either algorithms or randomization.

The Problem

When we studied NP-completeness, a core problem was 3-SAT: Given a set of clauses C_1, \dots, C_k , each of length 3, over a set of variables $X = \{x_1, \dots, x_n\}$, does there exist a satisfying truth assignment?

Intuitively, we can imagine such a problem arising in a system that tries to decide the truth or falsehood of statements about the world (the variables $\{x_i\}$), given pieces of information that relate them to one another (the clauses $\{C_j\}$). Now the world is a fairly contradictory place, and if our system gathers

enough information, it could well end up with a set of clauses that has no satisfying truth assignment. What then?

A natural approach, if we can't find a truth assignment that satisfies all clauses, is to turn the 3-SAT instance into an optimization problem: Given the set of input clauses C_1, \dots, C_k , find a truth assignment that satisfies *as many as possible*. We'll call this the *Maximum 3-Satisfiability Problem* (or *MAX 3-SAT* for short). Of course, this is an NP-hard optimization problem, since it's NP-complete to decide whether the maximum number of simultaneously satisfiable clauses is equal to k . Let's see what can be said about polynomial-time approximation algorithms.

Designing and Analyzing the Algorithm

A remarkably simple randomized algorithm turns out to give a strong performance guarantee for this problem. Suppose we set each variable x_1, \dots, x_n independently to 0 or 1 with probability $\frac{1}{2}$ each. What is the expected number of clauses satisfied by such a random assignment?

Let Z denote the random variable equal to the number of satisfied clauses. As in Section 13.3, let's decompose Z into a sum of random variables that each take the value 0 or 1; specifically, let $Z_i = 1$ if the clause C_i is satisfied, and 0 otherwise. Thus $Z = Z_1 + Z_2 + \dots + Z_k$. Now $E[Z_i]$ is equal to the probability that C_i is satisfied, and this can be computed easily as follows. In order for C_i *not* to be satisfied, each of its three variables must be assigned the value that fails to make it true; since the variables are set independently, the probability of this is $(\frac{1}{2})^3 = \frac{1}{8}$. Thus clause C_i is satisfied with probability $1 - \frac{1}{8} = \frac{7}{8}$, and so $E[Z_i] = \frac{7}{8}$.

Using linearity of expectation, we see that the expected number of satisfied clauses is $E[Z] = E[Z_1] + E[Z_2] + \dots + E[Z_k] = \frac{7}{8}k$. Since no assignment can satisfy more than k clauses, we have the following guarantee.

(13.14) *Consider a 3-SAT formula, where each clause has three different variables. The expected number of clauses satisfied by a random assignment is within an approximation factor $\frac{7}{8}$ of optimal.*

But, if we look at what really happened in the (admittedly simple) analysis of the random assignment, it's clear that something stronger is going on. For any random variable, there must be some point at which it assumes some value at least as large as its expectation. We've shown that for every instance of 3-SAT, a random truth assignment satisfies a $\frac{7}{8}$ fraction of all clauses in expectation; so, in particular, there must *exist* a truth assignment that satisfies a number of clauses that is at least as large as this expectation.

(13.15) *For every instance of 3-SAT, there is a truth assignment that satisfies at least a $\frac{7}{8}$ fraction of all clauses.*

There is something genuinely surprising about the statement of (13.15). We have arrived at a nonobvious fact about 3-SAT—the existence of an assignment satisfying many clauses—whose statement has nothing to do with randomization; but we have done so by a randomized construction. And, in fact, the randomized construction provides what is quite possibly the simplest proof of (13.15). This is a fairly widespread principle in the area of combinatorics—namely, that one can show the existence of some structure by showing that a random construction produces it with positive probability. Constructions of this sort are said to be applications of the *probabilistic method*.

Here’s a cute but minor application of (13.15): Every instance of 3-SAT with at most seven clauses is satisfiable. Why? If the instance has $k \leq 7$ clauses, then (13.15) implies that there is an assignment satisfying at least $\frac{7}{8}k$ of them. But when $k \leq 7$, it follows that $\frac{7}{8}k > k - 1$; and since the number of clauses satisfied by this assignment must be an integer, it must be equal to k . In other words, all clauses are satisfied.

Further Analysis: Waiting to Find a Good Assignment

Suppose we aren’t satisfied with a “one-shot” algorithm that produces a single assignment with a large number of satisfied clauses in expectation. Rather, we’d like a randomized algorithm whose expected running time is polynomial and that is guaranteed to output a truth assignment satisfying at least a $\frac{7}{8}$ fraction of all clauses.

A simple way to do this is to generate random truth assignments until one of them satisfies at least $\frac{7}{8}k$ clauses. We know that such an assignment exists, by (13.15); but how long will it take until we find one by random trials?

This is a natural place to apply the waiting-time bound we derived in (13.7). If we can show that the probability a random assignment satisfies at least $\frac{7}{8}k$ clauses is at least p , then the expected number of trials performed by the algorithm is $1/p$. So, in particular, we’d like to show that this quantity p is at least as large as an inverse polynomial in n and k .

For $j = 0, 1, 2, \dots, k$, let p_j denote the probability that a random assignment satisfies exactly j clauses. So the expected number of clauses satisfied, by the definition of expectation, is equal to $\sum_{j=0}^k j p_j$; and by the previous analysis, this is equal to $\frac{7}{8}k$. We are interested in the quantity $p = \sum_{j \geq 7k/8} p_j$. How can we use the lower bound on the expected value to give a lower bound on this quantity?

We start by writing

$$\frac{7}{8}k = \sum_{j=0}^k jp_j = \sum_{j < 7k/8} jp_j + \sum_{j \geq 7k/8} jp_j.$$

Now let k' denote the largest natural number that is strictly smaller than $\frac{7}{8}k$. The right-hand side of the above equation only increases if we replace the terms in the first sum by $k'p_j$ and the terms in the second sum by kp_j . We also observe that $\sum_{j < 7k/8} p_j = 1 - p$, and so

$$\frac{7}{8}k \leq \sum_{j < 7k/8} k'p_j + \sum_{j \geq 7k/8} kp_j = k'(1 - p) + kp \leq k' + kp,$$

and hence $kp \geq \frac{7}{8}k - k'$. But $\frac{7}{8}k - k' \geq \frac{1}{8}$, since k' is a natural number strictly smaller than $\frac{7}{8}$ times another natural number, and so

$$p \geq \frac{\frac{7}{8}k - k'}{k} \geq \frac{1}{8k}.$$

This was our goal—to get a lower bound on p —and so by the waiting-time bound (13.7), we see that the expected number of trials needed to find the satisfying assignment we want is at most $8k$.

(13.16) *There is a randomized algorithm with polynomial expected running time that is guaranteed to produce a truth assignment satisfying at least a $\frac{7}{8}$ fraction of all clauses.*

13.5 Randomized Divide and Conquer: Median-Finding and Quicksort

We've seen the divide-and-conquer paradigm for designing algorithms at various earlier points in the book. Divide and conquer often works well in conjunction with randomization, and we illustrate this by giving divide-and-conquer algorithms for two fundamental problems: computing the median of n numbers, and sorting. In each case, the “divide” step is performed using randomization; consequently, we will use expectations of random variables to analyze the time spent on recursive calls.



The Problem: Finding the Median

Suppose we are given a set of n numbers $S = \{a_1, a_2, \dots, a_n\}$. Their *median* is the number that would be in the middle position if we were to sort them. There's an annoying technical difficulty if n is even, since then there is no

“middle position”; thus we define things precisely as follows: The median of $S = \{a_1, a_2, \dots, a_n\}$ is equal to the k^{th} largest element in S , where $k = (n + 1)/2$ if n is odd, and $k = n/2$ if n is even. In what follows, we’ll assume for the sake of simplicity that all the numbers are distinct. Without this assumption, the problem becomes notationally more complicated, but no new ideas are brought into play.

It is clearly easy to compute the median in time $O(n \log n)$ if we simply sort the numbers first. But if one begins thinking about the problem, it’s far from clear why sorting is *necessary* for computing the median, or even why $\Omega(n \log n)$ time is necessary. In fact, we’ll show how a simple randomized approach, based on divide-and-conquer, yields an expected running time of $O(n)$.

Designing the Algorithm

A Generic Algorithm Based on Splitters The first key step toward getting an expected linear running time is to move from median-finding to the more general problem of *selection*. Given a set of n numbers S and a number k between 1 and n , consider the function $\text{Select}(S, k)$ that returns the k^{th} largest element in S . As special cases, Select includes the problem of finding the median of S via $\text{Select}(S, n/2)$ or $\text{Select}(S, (n + 1)/2)$; it also includes the easier problems of finding the minimum ($\text{Select}(S, 1)$) and the maximum ($\text{Select}(S, n)$). Our goal is to design an algorithm that implements Select so that it runs in expected time $O(n)$.

The basic structure of the algorithm implementing Select is as follows. We choose an element $a_i \in S$, the “splitter,” and form the sets $S^- = \{a_j : a_j < a_i\}$ and $S^+ = \{a_j : a_j > a_i\}$. We can then determine which of S^- or S^+ contains the k^{th} largest element, and iterate only on this one. Without specifying yet how we plan to choose the splitter, here’s a more concrete description of how we form the two sets and iterate.

```

Select( $S, k$ ):
  Choose a splitter  $a_i \in S$ 
  For each element  $a_j$  of  $S$ 
    Put  $a_j$  in  $S^-$  if  $a_j < a_i$ 
    Put  $a_j$  in  $S^+$  if  $a_j > a_i$ 
  Endfor
  If  $|S^-| = k - 1$  then
    The splitter  $a_i$  was in fact the desired answer
  Else if  $|S^-| \geq k$  then
    The  $k^{\text{th}}$  largest element lies in  $S^-$ 
    Recursively call  $\text{Select}(S^-, k)$ 

```

```

Else suppose  $|S^-| = \ell < k - 1$ 
    The  $k^{\text{th}}$  largest element lies in  $S^+$ 
    Recursively call Select( $S^+, k - 1 - \ell$ )
Endif

```

Observe that the algorithm is always called recursively on a strictly smaller set, so it must terminate. Also, observe that if $|S| = 1$, then we must have $k = 1$, and indeed the single element in S will be returned by the algorithm. Finally, from the choice of which recursive call to make, it's clear by induction that the right answer will be returned when $|S| > 1$ as well. Thus we have the following

(13.17) *Regardless of how the splitter is chosen, the algorithm above returns the k^{th} largest element of S .*

Choosing a Good Splitter Now let's consider how the running time of *Select* depends on the way we choose the splitter. Assuming we can select a splitter in linear time, the rest of the algorithm takes linear time plus the time for the recursive call. But how is the running time of the recursive call affected by the choice of the splitter? Essentially, it's important that the splitter significantly reduce the size of the set being considered, so that we don't keep making passes through large sets of numbers many times. So a good choice of splitter should produce sets S^- and S^+ that are approximately equal in size.

For example, if we could always choose the median as the splitter, then we could show a linear bound on the running time as follows. Let cn be the running time for *Select*, not counting the time for the recursive call. Then, with medians as splitters, the running time $T(n)$ would be bounded by the recurrence $T(n) \leq T(n/2) + cn$. This is a recurrence that we encountered at the beginning of Chapter 5, where we showed that it has the solution $T(n) = O(n)$.

Of course, hoping to be able to use the median as the splitter is rather circular, since the median is what we want to compute in the first place! But, in fact, one can show that any “well-centered” element can serve as a good splitter: If we had a way to choose splitters a_i such that there were at least εn elements both larger and smaller than a_i , for any fixed constant $\varepsilon > 0$, then the size of the sets in the recursive call would shrink by a factor of at least $(1 - \varepsilon)$ each time. Thus the running time $T(n)$ would be bounded by the recurrence $T(n) \leq T((1 - \varepsilon)n) + cn$. The same argument that showed the previous recurrence had the solution $T(n) = O(n)$ can be used here: If we unroll this recurrence for any $\varepsilon > 0$, we get

$$T(n) \leq cn + (1 - \varepsilon)cn + (1 - \varepsilon)^2 cn + \cdots = \left[1 + (1 - \varepsilon) + (1 - \varepsilon)^2 + \cdots \right] cn \leq \frac{1}{\varepsilon} \cdot cn,$$

since we have a convergent geometric series.

Indeed, the only thing to really beware of is a very “off-center” splitter. For example, if we always chose the minimum element as the splitter, then we may end up with a set in the recursive call that’s only one element smaller than we had before. In this case, the running time $T(n)$ would be bounded by the recurrence $T(n) \leq T(n - 1) + cn$. Unrolling this recurrence, we see that there’s a problem:

$$T(n) \leq cn + c(n - 1) + c(n - 2) + \cdots = \frac{cn(n + 1)}{2} = \Theta(n^2).$$

Random Splitters Choosing a “well-centered” splitter, in the sense we have just defined, is certainly similar in flavor to our original problem of choosing the median; but the situation is really not so bad, since *any* well-centered splitter will do.

Thus we will implement the as-yet-unspecified step of selecting a splitter using the following simple rule:

Choose a splitter $a_i \in S$ uniformly at random

The intuition here is very natural: since a fairly large fraction of the elements are reasonably well-centered, we will be likely to end up with a good splitter simply by choosing an element at random.

The analysis of the running time with a random splitter is based on this idea; we expect the size of the set under consideration to go down by a fixed constant fraction every iteration, so we should get a convergent series and hence a linear bound as previously. We now show how to make this precise.

Analyzing the Algorithm

We’ll say that the algorithm is in *phase* j when the size of the set under consideration is at most $n(\frac{3}{4})^j$ but greater than $n(\frac{3}{4})^{j+1}$. Let’s try to bound the expected time spent by the algorithm in phase j . In a given iteration of the algorithm, we say that an element of the set under consideration is *central* if at least a quarter of the elements are smaller than it and at least a quarter of the elements are larger than it.

Now observe that if a central element is chosen as a splitter, then at least a quarter of the set will be thrown away, the set will shrink by a factor of $\frac{3}{4}$ or better, and the current phase will come to an end. Moreover, half of all the

elements in the set are central, and so the probability that our random choice of splitter produces a central element is $\frac{1}{2}$. Hence, by our simple waiting-time bound (13.7), the expected number of iterations before a central element is found is 2; and so the expected number of iterations spent in phase j , for any j , is at most 2.

This is pretty much all we need for the analysis. Let X be a random variable equal to the number of steps taken by the algorithm. We can write it as the sum $X = X_0 + X_1 + X_2 + \dots$, where X_j is the expected number of steps spent by the algorithm in phase j . When the algorithm is in phase j , the set has size at most $n(\frac{3}{4})^j$, and so the number of steps required for one iteration in phase j is at most $cn(\frac{3}{4})^j$ for some constant c . We have just argued that the expected number of iterations spent in phase j is at most two, and hence we have $E[X_j] \leq 2cn(\frac{3}{4})^j$. Thus we can bound the total expected running time using linearity of expectation,

$$E[X] = \sum_j E[X_j] \leq \sum_j 2cn \left(\frac{3}{4}\right)^j = 2cn \sum_j \left(\frac{3}{4}\right)^j \leq 8cn,$$

since the sum $\sum_j (\frac{3}{4})^j$ is a geometric series that converges. Thus we have the following desired result.

(13.18) *The expected running time of `Select`(n, k) is $O(n)$.*

A Second Application: Quicksort

The randomized divide-and-conquer technique we used to find the median is also the basis of the sorting algorithm `Quicksort`. As before, we choose a splitter for the input set S , and separate S into the elements below the splitter value and those above it. The difference is that, rather than looking for the median on just one side of the splitter, we sort both sides recursively and glue the two sorted pieces together (with the splitter in between) to produce the overall output. Also, we need to explicitly include a base case for the recursive code: we only use recursion on sets of size at least 4. A complete description of `Quicksort` is as follows.

```

Quicksort( $S$ ):
  If  $|S| \leq 3$  then
    Sort  $S$ 
    Output the sorted list
  Else
    Choose a splitter  $a_i \in S$  uniformly at random
    For each element  $a_j$  of  $S$ 
```

```

        Put  $a_j$  in  $S^-$  if  $a_j < a_i$ 
        Put  $a_j$  in  $S^+$  if  $a_j > a_i$ 
    Endfor
    Recursively call Quicksort( $S^-$ ) and Quicksort( $S^+$ )
    Output the sorted set  $S^-$ , then  $a_i$ , then the sorted set  $S^+$ 
Endif

```

As with median-finding, the worst-case running time of this method is not so good. If we always select the smallest element as a splitter, then the running time $T(n)$ on n -element sets satisfies the same recurrence as before: $T(n) \leq T(n-1) + cn$, and so we end up with a time bound of $T(n) = \Theta(n^2)$. In fact, this is the worst-case running time for Quicksort.

On the positive side, if the splitters selected happened to be the medians of the sets at each iteration, then we get the recurrence $T(n) \leq 2T(n/2) + cn$, which arose frequently in the divide-and-conquer analyses of Chapter 5; the running time in this lucky case is $O(n \log n)$.

Here we are concerned with the *expected running time*; we will show that this can be bounded by $O(n \log n)$, almost as good as in the best case when the splitters are perfectly centered. Our analysis of Quicksort will closely follow the analysis of median-finding. Just as in the `Select` procedure that we used for median-finding, the crucial definition is that of a *central splitter*—one that divides the set so that each side contains at least a quarter of the elements. (As we discussed earlier, it is enough for the analysis that each side contains at least some fixed constant fraction of the elements; the use of a quarter here is chosen for convenience.) The idea is that a random choice is likely to lead to a central splitter, and central splitters work well. In the case of sorting, a central splitter divides the problem into two considerably smaller subproblems.

To simplify the presentation, we will slightly modify the algorithm so that it only issues its recursive calls when it finds a central splitter. Essentially, this modified algorithm differs from Quicksort in that it prefers to throw away an “off-center” splitter and try again; Quicksort, by contrast, launches the recursive calls even with an off-center splitter, and at least benefits from the work already done in splitting S . The point is that the expected running time of this modified algorithm can be analyzed very simply, by direct analogy with our analysis for median-finding. With a bit more work, a very similar but somewhat more involved analysis can also be done for the original Quicksort algorithm as well; however, we will not describe this analysis here.

```

Modified Quicksort( $S$ ):
    If  $|S| \leq 3$  then
        Sort  $S$ 
    
```

```

    Output the sorted list
Endif
Else
    While no central splitter has been found
        Choose a splitter  $a_i \in S$  uniformly at random
        For each element  $a_j$  of  $S$ 
            Put  $a_j$  in  $S^-$  if  $a_j < a_i$ 
            Put  $a_j$  in  $S^+$  if  $a_j > a_i$ 
        Endfor
        If  $|S^-| \geq |S|/4$  and  $|S^+| \geq |S|/4$  then
             $a_i$  is a central splitter
        Endif
    Endwhile
    Recursively call Quicksort( $S^-$ ) and Quicksort( $S^+$ )
    Output the sorted set  $S^-$ , then  $a_i$ , then the sorted set  $S^+$ 
Endif

```

Consider a subproblem for some set S . Each iteration of the **While** loop selects a possible splitter a_i and spends $O(|S|)$ time splitting the set and deciding if a_i is central. Earlier we argued that the number of iterations needed until we find a central splitter is at most 2. This gives us the following statement.

(13.19) *The expected running time for the algorithm on a set S , excluding the time spent on recursive calls, is $O(|S|)$.*

The algorithm is called recursively on multiple subproblems. We will group these subproblems by size. We'll say that the subproblem is of *type j* if the size of the set under consideration is at most $n(\frac{3}{4})^j$ but greater than $n(\frac{3}{4})^{j+1}$. By (13.19), the expected time spent on a subproblem of type j , excluding recursive calls, is $O(n(\frac{3}{4})^j)$. To bound the overall running time, we need to bound the number of subproblems for each type j . Splitting a type j subproblem via a central splitter creates two subproblems of higher type. So the subproblems of a given type j are disjoint. This gives us a bound on the number of subproblems.

(13.20) *The number of type j subproblems created by the algorithm is at most $(\frac{4}{3})^{j+1}$.*

There are at most $(\frac{4}{3})^{j+1}$ subproblems of type j , and the expected time spent on each is $O(n(\frac{3}{4})^j)$ by (13.19). Thus, by linearity of expectation, the expected time spent on subproblems of type j is $O(n)$. The number of different types is bounded by $\log_{\frac{4}{3}} n = O(\log n)$, which gives the desired bound.

(13.21) *The expected running time of Modified Quicksort is $O(n \log n)$.*

We considered this modified version of Quicksort to simplify the analysis. Coming back to the original Quicksort, our intuition suggests that the expected running time is no worse than in the modified algorithm, as accepting the noncentral splitters helps a bit with sorting, even if it does not help as much as when a central splitter is chosen. As mentioned earlier, one can in fact make this intuition precise, leading to an $O(n \log n)$ expected time bound for the original Quicksort algorithm; we will not go into the details of this here.

13.6 Hashing: A Randomized Implementation of Dictionaries

Randomization has also proved to be a powerful technique in the design of data structures. Here we discuss perhaps the most fundamental use of randomization in this setting, a technique called *hashing* that can be used to maintain a dynamically changing set of elements. In the next section, we will show how an application of this technique yields a very simple algorithm for a problem that we saw in Chapter 5—the problem of finding the closest pair of points in the plane.



The Problem

One of the most basic applications of data structures is to simply maintain a set of elements that changes over time. For example, such applications could include a large company maintaining the set of its current employees and contractors, a news indexing service recording the first paragraphs of news articles it has seen coming across the newswire, or a search algorithm keeping track of the small part of an exponentially large search space that it has already explored.

In all these examples, there is a *universe* U of possible elements that is extremely large: the set of all possible people, all possible paragraphs (say, up to some character length limit), or all possible solutions to a computationally hard problem. The data structure is trying to keep track of a set $S \subseteq U$ whose size is generally a negligible fraction of U , and the goal is to be able to insert and delete elements from S and quickly determine whether a given element belongs to S .

We will call a data structure that accomplishes this a *dictionary*. More precisely, a dictionary is a data structure that supports the following operations.

- **MakeDictionary.** This operation initializes a fresh dictionary that can maintain a subset S of U ; the dictionary starts out empty.
- **Insert(u)** adds element $u \in U$ to the set S . In many applications, there may be some additional information that we want to associate with u

(for example, u may be the name or ID number of an employee, and we want to also store some personal information about this employee), and we will simply imagine this being stored in the dictionary as part of a record together with u . (So, in general, when we talk about the element u , we really mean u and any additional information stored with u .)

- **Delete**(u) removes element u from the set S , if it is currently present.
- **Lookup**(u) determines whether u currently belongs to S ; if it does, it also retrieves any additional information stored with u .

Many of the implementations we've discussed earlier in the book involve (most of) these operations: For example, in the implementation of the BFS and DFS graph traversal algorithms, we needed to maintain the set S of nodes already visited. But there is a fundamental difference between those problems and the present setting, and that is the size of U . The universe U in BFS or DFS is the set of nodes V , which is already given explicitly as part of the input. Thus it is completely feasible in those cases to maintain a set $S \subseteq U$ as we did there: defining an array with $|U|$ positions, one for each possible element, and setting the array position for u equal to 1 if $u \in S$, and equal to 0 if $u \notin S$. This allows for insertion, deletion, and lookup of elements in constant time per operation, by simply accessing the desired array entry.

Here, by contrast, we are considering the setting in which the universe U is enormous. So we are not going to be able to use an array whose size is anywhere near that of U . The fundamental question is whether, in this case, we can still implement a dictionary to support the basic operations almost as quickly as when U was relatively small.

We now describe a randomized technique called *hashing* that addresses this question. While we will not be able to do quite as well as the case in which it is feasible to define an array over all of U , hashing will allow us to come quite close.

Designing the Data Structure

As a motivating example, let's think a bit more about the problem faced by an automated service that processes breaking news. Suppose you're receiving a steady stream of short articles from various wire services, weblog postings, and so forth, and you're storing the lead paragraph of each article (truncated to at most 1,000 characters). Because you're using many sources for the sake of full coverage, there's a lot of redundancy: the same article can show up many times.

When a new article shows up, you'd like to quickly check whether you've seen the lead paragraph before. So a dictionary is exactly what you want for this problem: The universe U is the set of all strings of length at most 1,000 (or of

length exactly 1,000, if we pad them out with blanks), and we're maintaining a set $S \subseteq U$ consisting of strings (i.e., lead paragraphs) that we've seen before.

One solution would be to keep a linked list of all paragraphs, and scan this list each time a new one arrives. But a Lookup operation in this case takes time proportional to $|S|$. How can we get back to something that looks like an array-based solution?

Hash Functions The basic idea of hashing is to work with an array of size $|S|$, rather than one comparable to the (astronomical) size of U .

Suppose we want to be able to store a set S of size up to n . We will set up an array H of size n to store the information, and use a function $h : U \rightarrow \{0, 1, \dots, n-1\}$ that maps elements of U to array positions. We call such a function h a *hash function*, and the array H a *hash table*. Now, if we want to add an element u to the set S , we simply place u in position $h(u)$ of the array H . In the case of storing paragraphs of text, we can think of $h(\cdot)$ as computing some kind of numerical signature or “check-sum” of the paragraph u , and this tells us the array position at which to store u .

This would work extremely well if, for all distinct u and v in our set S , it happened to be the case that $h(u) \neq h(v)$. In such a case, we could look up u in constant time: when we check array position $H[h(u)]$, it would either be empty or would contain just u .

In general, though, we cannot expect to be this lucky: there can be distinct elements $u, v \in S$ for which $h(u) = h(v)$. We will say that these two elements *collide*, since they are mapped to the same place in H . There are a number of ways to deal with collisions. Here we will assume that each position $H[i]$ of the hash table stores a linked list of all elements $u \in S$ with $h(u) = i$. The operation $\text{Lookup}(u)$ would now work as follows.

- Compute the hash function $h(u)$.
- Scan the linked list at position $H[h(u)]$ to see if u is present in this list.

Hence the time required for $\text{Lookup}(u)$ is proportional to the time to compute $h(u)$, plus the length of the linked list at $H[h(u)]$. And this latter quantity, in turn, is just the number of elements in S that collide with u . The **Insert** and **Delete** operations work similarly: **Insert** adds u to the linked list at position $H[h(u)]$, and **Delete** scans this list and removes u if it is present.

So now the goal is clear: We'd like to find a hash function that “spreads out” the elements being added, so that no one entry of the hash table H contains too many elements. This is not a problem for which worst-case analysis is very informative. Indeed, suppose that $|U| \geq n^2$ (we're imagining applications where it's much larger than this). Then, for any hash function h that we choose, there will be some set S of n elements that all map to the same

position. In the worst case, we will insert all the elements of this set, and then our Lookup operations will consist of scanning a linked list of length n .

Our main goal here is to show that randomization can help significantly for this problem. As usual, we won't make any assumptions about the set of elements S being random; we will simply exploit randomization in the design of the hash function. In doing this, we won't be able to completely avoid collisions, but can make them relatively rare enough, and so the lists will be quite short.

Choosing a Good Hash Function We've seen that the efficiency of the dictionary is based on the choice of the hash function h . Typically, we will think of U as a large set of numbers, and then use an easily computable function h that maps each number $u \in U$ to some value in the smaller range of integers $\{0, 1, \dots, n-1\}$. There are many simple ways to do this: we could use the first or last few digits of u , or simply take u modulo n . While these simple choices may work well in many situations, it is also possible to get large numbers of collisions. Indeed, a fixed choice of hash function may run into problems because of the types of elements u encountered in the application: Maybe the particular digits we use to define the hash function encode some property of u , and hence maybe only a few options are possible. Taking u modulo n can have the same problem, especially if n is a power of 2. To take a concrete example, suppose we used a hash function that took an English paragraph, used a standard character encoding scheme like ASCII to map it to a sequence of bits, and then kept only the first few bits in this sequence. We'd expect a huge number of collisions at the array entries corresponding to the bit strings that encoded common English words like *The*, while vast portions of the array can be occupied only by paragraphs that begin with strings like *qxf*, and hence will be empty.

A slightly better choice in practice is to take $(u \bmod p)$ for a prime number p that is approximately equal to n . While in some applications this may yield a good hashing function, it may not work well in all applications, and some primes may work much better than others (for example, primes very close to powers of 2 may not work so well).

Since hashing has been widely used in practice for a long time, there is a lot of experience with what makes for a good hash function, and many hash functions have been proposed that tend to work well empirically. Here we would like to develop a hashing scheme where we can prove that it results in efficient dictionary operations with high probability.

The basic idea, as suggested earlier, is to use randomization in the construction of h . First let's consider an extreme version of this: for every element $u \in U$, when we go to insert u into S , we select a value $h(u)$ uniformly at

random in the set $\{0, 1, \dots, n-1\}$, independently of all previous choices. In this case, the probability that two randomly selected values $h(u)$ and $h(v)$ are equal (and hence cause a collision) is quite small.

(13.22) *With this uniform random hashing scheme, the probability that two randomly selected values $h(u)$ and $h(v)$ collide—that is, that $h(u) = h(v)$ —is exactly $1/n$.*

Proof. Of the n^2 possible choices for the pair of values $(h(u), h(v))$, all are equally likely, and exactly n of these choices results in a collision. ■

However, it will not work to use a hash function with independently random chosen values. To see why, suppose we inserted u into S , and then later want to perform either `Delete(u)` or `Lookup(u)`. We immediately run into the “Where did I put it?” problem: We will need to know the random value $h(u)$ that we used, so we will need to have stored the value $h(u)$ in some form where we can quickly look it up. But this is exactly the same problem we were trying to solve in the first place.

There are two things that we can learn from (13.22). First, it provides a concrete basis for the intuition from practice that hash functions that spread things around in a “random” way can be effective at reducing collisions. Second, and more crucial for our goals here, we will be able to show how a more controlled use of randomization achieves performance as good as suggested in (13.22), but in a way that leads to an efficient dictionary implementation.

Universal Classes of Hash Functions The key idea is to choose a hash function at random not from the collection of all possible functions into $[0, n-1]$, but from a carefully selected class of functions. Each function h in our class of functions \mathcal{H} will map the universe U into the set $\{0, 1, \dots, n-1\}$, and we will design it so that it has two properties. First, we’d like it to come with the guarantee from (13.22):

- For any pair of elements $u, v \in U$, the probability that a randomly chosen $h \in \mathcal{H}$ satisfies $h(u) = h(v)$ is at most $1/n$.

We say that a class \mathcal{H} of functions is *universal* if it satisfies this first property. Thus (13.22) can be viewed as saying that the class of all possible functions from U into $\{0, 1, \dots, n-1\}$ is universal.

However, we also need \mathcal{H} to satisfy a second property. We will state this slightly informally for now and make it more precise later.

- Each $h \in \mathcal{H}$ can be compactly represented and, for a given $h \in \mathcal{H}$ and $u \in U$, we can compute the value $h(u)$ efficiently.

The class of all possible functions failed to have this property: Essentially, the only way to represent an arbitrary function from U into $\{0, 1, \dots, n-1\}$ is to write down the value it takes on every single element of U .

In the remainder of this section, we will show the surprising fact that there exist classes \mathcal{H} that satisfy both of these properties. Before we do this, we first make precise the basic property we need from a universal class of hash functions. We argue that if a function h is selected at random from a universal class of hash functions, then in any set $S \subset U$ of size at most n , and any $u \in U$, the expected number of items in S that collide with u is a constant.

(13.23) *Let \mathcal{H} be a universal class of hash functions mapping a universe U to the set $\{0, 1, \dots, n-1\}$, let S be an arbitrary subset of U of size at most n , and let u be any element in U . We define X to be a random variable equal to the number of elements $s \in S$ for which $h(s) = h(u)$, for a random choice of hash function $h \in \mathcal{H}$. (Here S and u are fixed, and the randomness is in the choice of $h \in \mathcal{H}$.) Then $E[X] \leq 1$.*

Proof. For an element $s \in S$, we define a random variable X_s that is equal to 1 if $h(s) = h(u)$, and equal to 0 otherwise. We have $E[X_s] = \Pr[X_s = 1] \leq 1/n$, since the class of functions is universal.

Now $X = \sum_{s \in S} X_s$, and so, by linearity of expectation, we have

$$E[X] = \sum_{s \in S} E[X_s] \leq |S| \cdot \frac{1}{n} \leq 1. \quad \blacksquare$$

Designing a Universal Class of Hash Functions Next we will design a universal class of hash functions. We will use a prime number $p \approx n$ as the size of the hash table H . To be able to use integer arithmetic in designing our hash functions, we will identify the universe with vectors of the form $x = (x_1, x_2, \dots, x_r)$ for some integer r , where $0 \leq x_i < p$ for each i . For example, we can first identify U with integers in the range $[0, N-1]$ for some N , and then use consecutive blocks of $\lfloor \log p \rfloor$ bits of u to define the corresponding coordinates x_i . If $U \subseteq [0, N-1]$, then we will need a number of coordinates $r \approx \log N / \log n$.

Let \mathcal{A} be the set of all vectors of the form $a = (a_1, \dots, a_r)$, where a_i is an integer in the range $[0, p-1]$ for each $i = 1, \dots, r$. For each $a \in \mathcal{A}$, we define the linear function

$$h_a(x) = \left(\sum_{i=1}^r a_i x_i \right) \bmod p.$$

This now completes our random implementation of dictionaries. We define the family of hash functions to be $\mathcal{H} = \{h_a : a \in \mathcal{A}\}$. To execute `MakeDictionary`, we choose a random hash function from \mathcal{H} ; in other words, we choose a random vector from \mathcal{A} (by choosing each coordinate uniformly at random), and form the function h_a . Note that in order to define \mathcal{A} , we need to find a prime number $p \geq n$. There are methods for generating prime numbers quickly, which we will not go into here. (In practice, this can also be accomplished using a table of known prime numbers, even for relatively large n .)

We then use this as the hash function with which to implement `Insert`, `Delete`, and `Lookup`. The family $\mathcal{H} = \{h_a : a \in \mathcal{A}\}$ satisfies a formal version of the second property we were seeking: It has a compact representation, since by simply choosing and remembering a random $a \in \mathcal{A}$, we can compute $h_a(u)$ for all elements $u \in U$. Thus, to show that \mathcal{H} leads to an efficient, hashing-based implementation of dictionaries, we just need to establish that \mathcal{H} is a universal family of hash functions.



Analyzing the Data Structure

If we are using a hash function h_a from the class \mathcal{H} that we've defined, then a collision $h_a(x) = h_a(y)$ defines a linear equation modulo the prime number p . In order to analyze such equations, it's useful to have the following "cancellation law."

(13.24) For any prime p and any integer $z \neq 0 \bmod p$, and any two integers α, β , if $\alpha z = \beta z \bmod p$, then $\alpha = \beta \bmod p$.

Proof. Suppose $\alpha z = \beta z \bmod p$. Then, by rearranging terms, we get $z(\alpha - \beta) = 0 \bmod p$, and hence $z(\alpha - \beta)$ is divisible by p . But $z \neq 0 \bmod p$, so z is not divisible by p . Since p is prime, it follows that $\alpha - \beta$ must be divisible by p ; that is, $\alpha = \beta \bmod p$ as claimed. ■

We now use this to prove the main result in our analysis.

(13.25) The class of linear functions \mathcal{H} defined above is universal.

Proof. Let $x = (x_1, x_2, \dots, x_r)$ and $y = (y_1, y_2, \dots, y_r)$ be two distinct elements of U . We need to show that the probability of $h_a(x) = h_a(y)$, for a randomly chosen $a \in \mathcal{A}$, is at most $1/p$.

Since $x \neq y$, then there must be an index j such that $x_j \neq y_j$. We now consider the following way of choosing the random vector $a \in \mathcal{A}$. We first choose all the coordinates a_i where $i \neq j$. Then, finally, we choose coordinate a_j . We will show that regardless of how all the other coordinates a_i were

chosen, the probability of $h_a(x) = h_a(y)$, taken over the final choice of a_j , is exactly $1/p$. It will follow that the probability of $h_a(x) = h_a(y)$ over the random choice of the full vector a must be $1/p$ as well.

This conclusion is intuitively clear: If the probability is $1/p$ regardless of how we choose all other a_i , then it is $1/p$ overall. There is also a direct proof of this using conditional probabilities. Let \mathcal{E} be the event that $h_a(x) = h_a(y)$, and let \mathcal{F}_b be the event that all coordinates a_i (for $i \neq j$) receive a sequence of values b . We will show, below, that $\Pr[\mathcal{E} \mid \mathcal{F}_b] = 1/p$ for all b . It then follows that $\Pr[\mathcal{E}] = \sum_b \Pr[\mathcal{E} \mid \mathcal{F}_b] \cdot \Pr[\mathcal{F}_b] = (1/p) \sum_b \Pr[\mathcal{F}_b] = 1/p$.

So, to conclude the proof, we assume that values have been chosen arbitrarily for all other coordinates a_i , and we consider the probability of selecting a_j so that $h_a(x) = h_a(y)$. By rearranging terms, we see that $h_a(x) = h_a(y)$ if and only if

$$a_j(y_j - x_j) = \sum_{i \neq j} a_i(x_i - y_i) \bmod p.$$

Since the choices for all a_i ($i \neq j$) have been fixed, we can view the right-hand side as some fixed quantity m . Also, let us define $z = y_j - x_j$.

Now it is enough to show that there is exactly one value $0 \leq a_j < p$ that satisfies $a_j z = m \bmod p$; indeed, if this is the case, then there is a probability of exactly $1/p$ of choosing this value for a_j . So suppose there were two such values, a_j and a'_j . Then we would have $a_j z = a'_j z \bmod p$, and so by (13.24) we would have $a_j = a'_j \bmod p$. But we assumed that $a_j, a'_j < p$, and so in fact a_j and a'_j would be the same. It follows that there is only one a_j in this range that satisfies $a_j z = m \bmod p$.

Tracing back through the implications, this means that the probability of choosing a_j so that $h_a(x) = h_a(y)$ is $1/p$, however we set the other coordinates a_i in a ; thus the probability that x and y collide is $1/p$. Thus we have shown that \mathcal{H} is a universal class of hash functions. ■

13.7 Finding the Closest Pair of Points: A Randomized Approach

In Chapter 5, we used the divide-and-conquer technique to develop an $O(n \log n)$ time algorithm for the problem of finding the closest pair of points in the plane. Here we will show how to use randomization to develop a different algorithm for this problem, using an underlying dictionary data structure. We will show that this algorithm runs in $O(n)$ expected time, plus $O(n)$ expected dictionary operations.

There are several related reasons why it is useful to express the running time of our algorithm in this way, accounting for the dictionary operations

This is the probability of a miss on the request to s . Summing over all requests to unmarked items, we have

$$E[X_j] \leq c_j + \sum_{i=1}^{k-c_j} \frac{c_j}{k-i+1} \leq c_j \left[1 + \sum_{\ell=c_j+1}^k \frac{1}{\ell} \right] = c_j(1 + H(k) - H(c_j)) \leq c_j H(k).$$

Thus the total expected number of misses incurred by the Randomized Marking Algorithm is

$$E[M_\sigma] = \sum_{j=1}^r E[X_j] \leq H(k) \sum_{j=1}^r c_j. \quad \blacksquare$$

Combining (13.39) and (13.40), we immediately get the following performance guarantee.

(13.41) *The expected number of misses incurred by the Randomized Marking Algorithm is at most $2H(k) \cdot f(\sigma) = O(\log k) \cdot f(\sigma)$.*

13.9 Chernoff Bounds

In Section 13.3, we defined the expectation of a random variable formally and have worked with this definition and its consequences ever since. Intuitively, we have a sense that the value of a random variable ought to be “near” its expectation with reasonably high probability, but we have not yet explored the extent to which this is true. We now turn to some results that allow us to reach conclusions like this, and see a sampling of the applications that follow.

We say that two random variables X and Y are *independent* if, for any values i and j , the events $\Pr[X=i]$ and $\Pr[Y=j]$ are independent. This definition extends naturally to larger sets of random variables. Now consider a random variable X that is a sum of several independent 0-1-valued random variables: $X = X_1 + X_2 + \cdots + X_n$, where X_i takes the value 1 with probability p_i , and the value 0 otherwise. By linearity of expectation, we have $E[X] = \sum_{i=1}^n p_i$. Intuitively, the independence of the random variables X_1, X_2, \dots, X_n suggests that their fluctuations are likely to “cancel out,” and so their sum X will have a value close to its expectation with high probability. This is in fact true, and we state two concrete versions of this result: one bounding the probability that X deviates above $E[X]$, the other bounding the probability that X deviates below $E[X]$. We call these results *Chernoff bounds*, after one of the probabilists who first established bounds of this form.

(13.42) Let X, X_1, X_2, \dots, X_n be defined as above, and assume that $\mu \geq E[X]$. Then, for any $\delta > 0$, we have

$$\Pr[X > (1 + \delta)\mu] < \left[\frac{e^\delta}{(1 + \delta)^{(1 + \delta)}} \right]^\mu.$$

Proof. To bound the probability that X exceeds $(1 + \delta)\mu$, we go through a sequence of simple transformations. First note that, for any $t > 0$, we have $\Pr[X > (1 + \delta)\mu] = \Pr[e^{tX} > e^{t(1 + \delta)\mu}]$, as the function $f(x) = e^{tx}$ is monotone in x . We will use this observation with a t that we'll select later.

Next we use some simple properties of the expectation. For a random variable Y , we have $\gamma \Pr[Y > \gamma] \leq E[Y]$, by the definition of the expectation. This allows us to bound the probability that Y exceeds γ in terms of $E[Y]$. Combining these two ideas, we get the following inequalities.

$$\Pr[X > (1 + \delta)\mu] = \Pr[e^{tX} > e^{t(1 + \delta)\mu}] \leq e^{-t(1 + \delta)\mu} E[e^{tX}].$$

Next we need to bound the expectation $E[e^{tX}]$. Writing X as $X = \sum_i X_i$, the expectation is $E[e^{tX}] = E[e^{t \sum_i X_i}] = E[\prod_i e^{tX_i}]$. For independent variables Y and Z , the expectation of the product YZ is $E[YZ] = E[Y]E[Z]$. The variables X_i are independent, so we get $E[\prod_i e^{tX_i}] = \prod_i E[e^{tX_i}]$.

Now, e^{tX_i} is e^t with probability p_i and $e^0 = 1$ otherwise, so its expectation can be bounded as

$$E[e^{tX_i}] = p_i e^t + (1 - p_i) = 1 + p_i(e^t - 1) \leq e^{p_i(e^t - 1)},$$

where the last inequality follows from the fact that $1 + \alpha \leq e^\alpha$ for any $\alpha \geq 0$. Combining the inequalities, we get the following bound.

$$\begin{aligned} \Pr[X > (1 + \delta)\mu] &\leq e^{-t(1 + \delta)\mu} E[e^{tX}] = e^{-t(1 + \delta)\mu} \prod_i E[e^{tX_i}] \\ &\leq e^{-t(1 + \delta)\mu} \prod_i e^{p_i(e^t - 1)} \leq e^{-t(1 + \delta)\mu} e^{\mu(e^t - 1)}. \end{aligned}$$

To obtain the bound claimed by the statement, we substitute $t = \ln(1 + \delta)$. ■

Where (13.42) provided an upper bound, showing that X is not likely to deviate far above its expectation, the next statement, (13.43), provides a lower bound, showing that X is not likely to deviate far below its expectation. Note that the statements of the results are not symmetric, and this makes sense: For the upper bound, it is interesting to consider values of δ much larger than 1, while this would not make sense for the lower bound.

(13.43) Let X, X_1, X_2, \dots, X_n and μ be defined as above. Then for any $1 > \delta > 0$, we have

$$\Pr [X < (1 - \delta)\mu] < e^{-\frac{1}{2}\mu\delta^2}.$$

The proof of (13.43) is similar to the proof of (13.42), and we do not give it here. For the applications that follow, the statements of (13.42) and (13.43), rather than the internals of their proofs, are the key things to keep in mind.

13.10 Load Balancing

In Section 13.1, we considered a distributed system in which communication among processes was difficult, and randomization to some extent replaced explicit coordination and synchronization. We now revisit this theme through another stylized example of randomization in a distributed setting.



The Problem

Suppose we have a system in which m jobs arrive in a stream and need to be processed immediately. We have a collection of n identical processors that are capable of performing the jobs; so the goal is to assign each job to a processor in a way that balances the workload evenly across the processors. If we had a central controller for the system that could receive each job and hand it off to the processors in round-robin fashion, it would be trivial to make sure that each processor received at most $\lceil m/n \rceil$ jobs—the most even balancing possible.

But suppose the system lacks the coordination or centralization to implement this. A much more lightweight approach would be to simply assign each job to one of the processors uniformly at random. Intuitively, this should also balance the jobs evenly, since each processor is equally likely to get each job. At the same time, since the assignment is completely random, one doesn't expect everything to end up perfectly balanced. So we ask: How well does this simple randomized approach work?

Although we will stick to the motivation in terms of jobs and processors here, it is worth noting that comparable issues come up in the analysis of hash functions, as we saw in Section 13.6. There, instead of assigning jobs to processors, we're assigning elements to entries in a hash table. The concern about producing an even balancing in the case of hash tables is based on wanting to keep the number of collisions at any particular entry relatively small. As a result, the analysis in this section is also relevant to the study of hashing schemes.

- 2^{t-2} ways to get an odd sum on the first $t - 1$ variables (by induction), followed by an assignment of 1 to the t^{th} .

The remaining 2^{t-1} assignments give an odd sum, and this completes the induction step.

Once we have $E[X_r] = 1/2$, we conclude as in part (a): Linearity of expectation gives us $E[X] = \sum_r E[X_r] = k/2 \geq c^*/2$.

Exercises

1. *3-Coloring* is a yes/no question, but we can phrase it as an optimization problem as follows.

Suppose we are given a graph $G = (V, E)$, and we want to color each node with one of three colors, even if we aren't necessarily able to give different colors to every pair of adjacent nodes. Rather, we say that an edge (u, v) is *satisfied* if the colors assigned to u and v are different.

Consider a 3-coloring that maximizes the number of satisfied edges, and let c^* denote this number. Give a polynomial-time algorithm that produces a 3-coloring that satisfies at least $\frac{2}{3}c^*$ edges. If you want, your algorithm can be randomized; in this case, the *expected* number of edges it satisfies should be at least $\frac{2}{3}c^*$.

2. Consider a county in which 100,000 people vote in an election. There are only two candidates on the ballot: a Democratic candidate (denoted D) and a Republican candidate (denoted R). As it happens, this county is heavily Democratic, so 80,000 people go to the polls with the intention of voting for D , and 20,000 go to the polls with the intention of voting for R .

However, the layout of the ballot is a little confusing, so each voter, independently and with probability $\frac{1}{100}$, votes for the wrong candidate—that is, the one that he or she *didn't* intend to vote for. (Remember that in this election, there are only two candidates on the ballot.)

Let X denote the random variable equal to the number of votes received by the Democratic candidate D , when the voting is conducted with this process of error. Determine the expected value of X , and give an explanation of your derivation of this value.

3. In Section 13.1, we saw a simple distributed protocol to solve a particular contention-resolution problem. Here is another setting in which randomization can help with contention resolution, through the distributed construction of an independent set.

Suppose we have a system with n processes. Certain pairs of processes are in *conflict*, meaning that they both require access to a shared resource. In a given time interval, the goal is to schedule a large subset S of the processes to run—the rest will remain idle—so that no two conflicting processes are both in the scheduled set S . We'll call such a set S *conflict-free*.

One can picture this process in terms of a graph $G = (V, E)$ with a node representing each process and an edge joining pairs of processes that are in conflict. It is easy to check that a set of processes S is conflict-free if and only if it forms an independent set in G . This suggests that finding a maximum-size conflict-free set S , for an arbitrary conflict G , will be difficult (since the general Independent Set Problem is reducible to this problem). Nevertheless, we can still look for heuristics that find a reasonably large conflict-free set. Moreover, we'd like a simple method for achieving this without centralized control: Each process should communicate with only a small number of other processes and then decide whether or not it should belong to the set S .

We will suppose for purposes of this question that each node has exactly d neighbors in the graph G . (That is, each process is in conflict with exactly d other processes.)

- (a) Consider the following simple protocol.

Each process P_i independently picks a random value x_i ; it sets x_i to 1 with probability $\frac{1}{2}$ and sets x_i to 0 with probability $\frac{1}{2}$. It then decides to enter the set S if and only if it chooses the value 1, and each of the processes with which it is in conflict chooses the value 0.

Prove that the set S resulting from the execution of this protocol is conflict-free. Also, give a formula for the expected size of S in terms of n (the number of processes) and d (the number of conflicts per process).

- (b) The choice of the probability $\frac{1}{2}$ in the protocol above was fairly arbitrary, and it's not clear that it should give the best system performance. A more general specification of the protocol would replace the probability $\frac{1}{2}$ by a parameter p between 0 and 1, as follows.

Each process P_i independently picks a random value x_i ; it sets x_i to 1 with probability p and sets x_i to 0 with probability $1 - p$. It then decides to enter the set S if and only if it chooses the value 1, and each of the processes with which it is in conflict chooses the value 0.

In terms of the parameters of the graph G , give a value of p so that the expected size of the resulting set S is as large as possible. Give a formula for the expected size of S when p is set to this optimal value.

4. A number of *peer-to-peer systems* on the Internet are based on *overlay networks*. Rather than using the physical Internet topology as the network on which to perform computation, these systems run protocols by which nodes choose collections of virtual “neighbors” so as to define a higher-level graph whose structure may bear little or no relation to the underlying physical network. Such an overlay network is then used for sharing data and services, and it can be extremely flexible compared with a physical network, which is hard to modify in real time to adapt to changing conditions.

Peer-to-peer networks tend to grow through the arrival of new participants, who join by linking into the existing structure. This growth process has an intrinsic effect on the characteristics of the overall network. Recently, people have investigated simple abstract models for network growth that might provide insight into the way such processes behave, at a qualitative level, in real networks.

Here’s a simple example of such a model. The system begins with a single node v_1 . Nodes then join one at a time; as each node joins, it executes a protocol whereby it forms a directed link to a single other node chosen uniformly at random from those already in the system. More concretely, if the system already contains nodes v_1, v_2, \dots, v_{k-1} and node v_k wishes to join, it randomly selects one of v_1, v_2, \dots, v_{k-1} and links to this node.

Suppose we run this process until we have a system consisting of nodes v_1, v_2, \dots, v_n ; the random process described above will produce a directed network in which each node other than v_1 has exactly one outgoing edge. On the other hand, a node may have multiple incoming links, or none at all. The incoming links to a node v_j reflect all the other nodes whose access into the system is via v_j ; so if v_j has many incoming links, this can place a large load on it. To keep the system load-balanced, then, we’d like all nodes to have a roughly comparable number of incoming links. That’s unlikely to happen here, however, since nodes that join earlier in the process are likely to have more incoming links than nodes that join later. Let’s try to quantify this imbalance as follows.

- (a) Given the random process described above, what is the expected number of incoming links to node v_j in the resulting network? Give an exact formula in terms of n and j , and also try to express this quantity

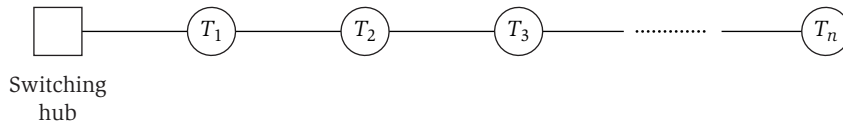


Figure 13.6 Towns T_1, T_2, \dots, T_n need to decide how to share the cost of the cable.

asymptotically (via an expression without large summations) using $\Theta(\cdot)$ notation.

- (b) Part (a) makes precise a sense in which the nodes that arrive early carry an “unfair” share of the connections in the network. Another way to quantify the imbalance is to observe that, in a run of this random process, we expect many nodes to end up with no incoming links.

Give a formula for the expected number of nodes with no incoming links in a network grown randomly according to this model.

5. Out in a rural part of the county somewhere, n small towns have decided to get connected to a large Internet switching hub via a high-volume fiber-optic cable. The towns are labeled T_1, T_2, \dots, T_n , and they are all arranged on a single long highway, so that town T_i is i miles from the switching hub (See Figure 13.6).

Now this cable is quite expensive; it costs k dollars per mile, resulting in an overall cost of kn dollars for the whole cable. The towns get together and discuss how to divide up the cost of the cable.

First, one of the towns way out at the far end of the highway makes the following proposal.

Proposal A. *Divide the cost evenly among all towns, so each pays k dollars.*

There’s some sense in which Proposal A is fair, since it’s as if each town is paying for the mile of cable directly leading up to it.

But one of the towns very close to the switching hub objects, pointing out that the faraway towns are actually benefiting from a large section of the cable, whereas the close-in towns only benefit from a short section of it. So they make the following counterproposal.

Proposal B. *Divide the cost so that the contribution of town T_i is proportional to i , its distance from the switching hub.*

One of the other towns very close to the switching hub points out that there’s another way to do a nonproportional division that is also

natural. This is based on conceptually dividing the cable into n equal-length “edges” e_1, \dots, e_n , where the first edge e_1 runs from the switching hub to T_1 , and the i^{th} edge e_i ($i > 1$) runs from T_{i-1} to T_i . Now we observe that, while all the towns benefit from e_1 , only the last town benefits from e_n . So they suggest

Proposal C. *Divide the cost separately for each edge e_i . The cost of e_i should be shared equally by the towns T_i, T_{i+1}, \dots, T_n , since these are the towns “downstream” of e_i .*

So now the towns have many different options; which is the fairest? To resolve this, they turn to the work of Lloyd Shapley, one of the most famous mathematical economists of the 20th century. He proposed what is now called the *Shapley value* as a general mechanism for sharing costs or benefits among several parties. It can be viewed as determining the “marginal contribution” of each party, *assuming the parties arrive in a random order*.

Here’s how it would work concretely in our setting. Consider an ordering \mathcal{O} of the towns, and suppose that the towns “arrive” in this order. The *marginal cost of town T_i in order \mathcal{O}* is determined as follows. If T_i is first in the order \mathcal{O} , then T_i pays ki , the cost of running the cable all the way from the switching hub to T_i . Otherwise, look at the set of towns that come before T_i in the order \mathcal{O} , and let T_j be the farthest among these towns from the switching hub. When T_i arrives, we assume the cable already reaches out to T_j but no farther. So if $j > i$ (T_j is farther out than T_i), then the marginal cost of T_i is 0, since the cable already runs past T_i on its way out to T_j . On the other hand, if $j < i$, then the marginal cost of T_i is $k(i - j)$: the cost of extending the cable from T_j out to T_i .

(For example, suppose $n = 3$ and the towns arrive in the order T_1, T_3, T_2 . First T_1 pays k when it arrives. Then, when T_3 arrives, it only has to pay $2k$ to extend the cable from T_1 . Finally, when T_2 arrives, it doesn’t have to pay anything since the cable already runs past it out to T_3 .)

Now, let X_i be the random variable equal to the marginal cost of town T_i when the order \mathcal{O} is selected uniformly at random from all permutations of the towns. Under the rules of the Shapley value, the amount that T_i should contribute to the overall cost of the cable is the expected value of X_i .

The question is: Which of the three proposals above, if any, gives the same division of costs as the Shapley value cost-sharing mechanism? Give a proof for your answer.

6. One of the (many) hard problems that arises in genome mapping can be formulated in the following abstract way. We are given a set of n *markers* $\{\mu_1, \dots, \mu_n\}$ —these are positions on a chromosome that we are trying to map—and our goal is to output a linear ordering of these markers. The output should be consistent with a set of k *constraints*, each specified by a triple (μ_i, μ_j, μ_k) , requiring that μ_j lie *between* μ_i and μ_k in the total ordering that we produce. (Note that this constraint does not specify which of μ_i or μ_k should come first in the ordering, only that μ_j should come between them.)

Now it is not always possible to satisfy all constraints simultaneously, so we wish to produce an ordering that satisfies as many as possible. Unfortunately, deciding whether there is an ordering that satisfies at least k' of the k constraints is an NP-complete problem (you don't have to prove this.)

Give a constant $\alpha > 0$ (independent of n) and an algorithm with the following property. If it is possible to satisfy k^* of the constraints, then the algorithm produces an ordering of markers satisfying at least αk^* of the constraints. Your algorithm may be randomized; in this case it should produce an ordering for which the *expected* number of satisfied constraints is at least αk^* .

7. In Section 13.4, we designed an approximation algorithm to within a factor of $7/8$ for the MAX 3-SAT Problem, where we assumed that each clause has terms associated with three different variables. In this problem, we will consider the analogous MAX SAT Problem: Given a set of clauses C_1, \dots, C_k over a set of variables $X = \{x_1, \dots, x_n\}$, find a truth assignment satisfying as many of the clauses as possible. Each clause has at least one term in it, and all the variables in a single clause are distinct, but otherwise we do not make any assumptions on the length of the clauses: There may be clauses that have a lot of variables, and others may have just a single variable.
- First consider the randomized approximation algorithm we used for MAX 3-SAT, setting each variable independently to *true* or *false* with probability $1/2$ each. Show that the expected number of clauses satisfied by this random assignment is at least $k/2$, that is, at least half of the clauses are satisfied in expectation. Give an example to show that there are MAX SAT instances such that no assignment satisfies more than half of the clauses.
 - If we have a clause that consists only of a single term (e.g., a clause consisting just of x_1 , or just of \bar{x}_2), then there is only a single way to satisfy it: We need to set the corresponding variable in the appropriate

way. If we have two clauses such that one consists of just the term x_i , and the other consists of just the negated term \bar{x}_i , then this is a pretty direct contradiction.

Assume that our instance has no such pair of “conflicting clauses”; that is, for no variable x_i do we have both a clause $C = \{x_i\}$ and a clause $C' = \{\bar{x}_i\}$. Modify the randomized procedure above to improve the approximation factor from $1/2$ to at least $.6$. That is, change the algorithm so that the expected number of clauses satisfied by the process is at least $.6k$.

- (c) Give a randomized polynomial-time algorithm for the general MAX SAT Problem, so that the expected number of clauses satisfied by the algorithm is at least a $.6$ fraction of the maximum possible.

(Note that, by the example in part (a), there are instances where one cannot satisfy more than $k/2$ clauses; the point here is that we’d still like an efficient algorithm that, in expectation, can satisfy a $.6$ fraction of the maximum that can be satisfied by an optimal assignment.)

8. Let $G = (V, E)$ be an undirected graph with n nodes and m edges. For a subset $X \subseteq V$, we use $G[X]$ to denote the subgraph *induced* on X —that is, the graph whose node set is X and whose edge set consists of all edges of G for which both ends lie in X .

We are given a natural number $k \leq n$ and are interested in finding a set of k nodes that induces a “dense” subgraph of G ; we’ll phrase this concretely as follows. Give a polynomial-time algorithm that produces, for a given natural number $k \leq n$, a set $X \subseteq V$ of k nodes with the property that the induced subgraph $G[X]$ has at least $\frac{mk(k-1)}{n(n-1)}$ edges.

You may give either (a) a deterministic algorithm, or (b) a randomized algorithm that has an expected running time that is polynomial, and that only outputs correct answers.

9. Suppose you’re designing strategies for selling items on a popular auction Web site. Unlike other auction sites, this one uses a *one-pass auction*, in which each bid must be immediately (and irrevocably) accepted or refused. Specifically, the site works as follows.
- First a seller puts up an item for sale.
 - Then buyers appear in sequence.
 - When buyer i appears, he or she makes a bid $b_i > 0$.
 - The seller must decide immediately whether to accept the bid or not. If the seller accepts the bid, the item is sold and all future buyers are

turned away. If the seller rejects the bid, buyer i departs and the bid is withdrawn; and only then does the seller see any future buyers.

Suppose an item is offered for sale, and there are n buyers, each with a distinct bid. Suppose further that the buyers appear in a random order, and that the seller knows the number n of buyers. We'd like to design a strategy whereby the seller has a reasonable chance of accepting the highest of the n bids. By a *strategy*, we mean a rule by which the seller decides whether to accept each presented bid, based only on the value of n and the sequence of bids seen so far.

For example, the seller could always accept the first bid presented. This results in the seller accepting the highest of the n bids with probability only $1/n$, since it requires the highest bid to be the first one presented.

Give a strategy under which the seller accepts the highest of the n bids with probability at least $1/4$, regardless of the value of n . (For simplicity, you may assume that n is an even number.) Prove that your strategy achieves this probabilistic guarantee.

10. Consider a very simple online auction system that works as follows. There are n *bidding agents*; agent i has a bid b_i , which is a positive natural number. We will assume that all bids b_i are distinct from one another. The bidding agents appear in an order chosen uniformly at random, each proposes its bid b_i in turn, and at all times the system maintains a variable b^* equal to the highest bid seen so far. (Initially b^* is set to 0.)

What is the expected number of times that b^* is updated when this process is executed, as a function of the parameters in the problem?

Example. Suppose $b_1 = 20$, $b_2 = 25$, and $b_3 = 10$, and the bidders arrive in the order 1, 3, 2. Then b^* is updated for 1 and 2, but not for 3.

11. *Load balancing algorithms* for parallel or distributed systems seek to spread out collections of computing jobs over multiple machines. In this way, no one machine becomes a "hot spot." If some kind of central coordination is possible, then the load can potentially be spread out almost perfectly. But what if the jobs are coming from diverse sources that can't coordinate? As we saw in Section 13.10, one option is to assign them to machines at random and hope that this randomization will work to prevent imbalances. Clearly, this won't generally work as well as a perfectly centralized solution, but it can be quite effective. Here we try analyzing some variations and extensions on the simple load balancing heuristic we considered in Section 13.10.

Suppose you have k machines, and k jobs show up for processing. Each job is assigned to one of the k machines independently at random (with each machine equally likely).

- (a) Let $N(k)$ be the expected number of machines that do not receive any jobs, so that $N(k)/k$ is the expected fraction of machines with nothing to do. What is the value of the limit $\lim_{k \rightarrow \infty} N(k)/k$? Give a proof of your answer.
 - (b) Suppose that machines are not able to queue up excess jobs, so if the random assignment of jobs to machines sends more than one job to a machine M , then M will do the first of the jobs it receives and reject the rest. Let $R(k)$ be the expected number of rejected jobs; so $R(k)/k$ is the expected fraction of rejected jobs. What is $\lim_{k \rightarrow \infty} R(k)/k$? Give a proof of your answer.
 - (c) Now assume that machines have slightly larger buffers; each machine M will do the first two jobs it receives, and reject any additional jobs. Let $R_2(k)$ denote the expected number of rejected jobs under this rule. What is $\lim_{k \rightarrow \infty} R_2(k)/k$? Give a proof of your answer.
12. Consider the following analogue of Karger's algorithm for finding minimum s - t cuts. We will contract edges iteratively using the following randomized procedure. In a given iteration, let s and t denote the possibly contracted nodes that contain the original nodes s and t , respectively. To make sure that s and t do not get contracted, at each iteration we delete any edges connecting s and t and select a random edge to contract among the remaining edges. Give an example to show that the probability that this method finds a minimum s - t cut can be exponentially small.
 13. Consider a balls-and-bins experiment with $2n$ balls but only two bins. As usual, each ball independently selects one of the two bins, both bins equally likely. The expected number of balls in each bin is n . In this problem, we explore the question of how big their difference is likely to be. Let X_1 and X_2 denote the number of balls in the two bins, respectively. (X_1 and X_2 are random variables.) Prove that for any $\varepsilon > 0$ there is a constant $c > 0$ such that the probability $\Pr[X_1 - X_2 > c\sqrt{n}] \leq \varepsilon$.
 14. Some people designing parallel physical simulations come to you with the following problem. They have a set P of k *basic processes* and want to assign each process to run on one of two machines, M_1 and M_2 . They are then going to run a sequence of n *jobs*, J_1, \dots, J_n . Each job J_i is represented by a set $P_i \subseteq P$ of exactly $2n$ basic processes which must be running (each on its assigned machine) while the job is processed. An assignment of basic processes to machines will be called *perfectly balanced* if, for

each job J_i , exactly n of the basic processes associated with J_i have been assigned to each of the two machines. An assignment of basic processes to machines will be called *nearly balanced* if, for each job J_i , no more than $\frac{4}{3}n$ of the basic processes associated with J_i have been assigned to the same machine.

- (a) Show that for arbitrarily large values of n , there exist sequences of jobs J_1, \dots, J_n for which no perfectly balanced assignment exists.
 - (b) Suppose that $n \geq 200$. Give an algorithm that takes an arbitrary sequence of jobs J_1, \dots, J_n and produces a nearly balanced assignment of basic processes to machines. Your algorithm may be randomized, in which case its expected running time should be polynomial, and it should always produce the correct answer.
15. Suppose you are presented with a very large set S of real numbers, and you'd like to approximate the median of these numbers by sampling. You may assume all the numbers in S are distinct. Let $n = |S|$; we will say that a number x is an ε -approximate median of S if at least $(\frac{1}{2} - \varepsilon)n$ numbers in S are less than x , and at least $(\frac{1}{2} - \varepsilon)n$ numbers in S are greater than x .
- Consider an algorithm that works as follows. You select a subset $S' \subseteq S$ uniformly at random, compute the median of S' , and return this as an approximate median of S . Show that there is an absolute constant c , independent of n , so that if you apply this algorithm with a sample S' of size c , then with probability at least .99, the number returned will be a (.05)-approximate median of S . (You may consider either the version of the algorithm that constructs S' by sampling with replacement, so that an element of S can be selected multiple times, or one without replacement.)
16. Consider the following (partially specified) method for transmitting a message securely between a sender and a receiver. The message will be represented as a string of bits. Let $\Sigma = \{0, 1\}$, and let Σ^* denote the set of all strings of 0 or more bits (e.g., $0, 00, 1110001 \in \Sigma^*$). The "empty string," with no bits, will be denoted $\lambda \in \Sigma^*$.

The sender and receiver share a secret function $f : \Sigma^* \times \Sigma \rightarrow \Sigma$. That is, f takes a word and a bit, and returns a bit. When the receiver gets a sequence of bits $\alpha \in \Sigma^*$, he or she runs the following method to decipher it.

Let $\alpha = \alpha_1\alpha_2 \dots \alpha_n$, where n is the number of bits in α

The goal is to produce an n -bit deciphered message,

$\beta = \beta_1\beta_2 \dots \beta_n$

Set $\beta_1 = f(\lambda, \alpha_1)$

```

For  $i = 2, 3, 4, \dots, n$ 
  Set  $\beta_i = f(\beta_1 \beta_2 \dots \beta_{i-1}, \alpha_i)$ 
Endfor
Output  $\beta$ 

```

One could view this as a type of “stream cipher with feedback.” One problem with this approach is that, if any bit α_i gets corrupted in transmission, it will corrupt the computed value of β_j for all $j \geq i$.

We consider the following problem. A sender S wants to transmit the same (plain-text) message β to each of k receivers R_1, \dots, R_k . With each one, he shares a different secret function $f^{(i)}$. Thus he sends a different encrypted message $\alpha^{(i)}$ to each receiver, so that $\alpha^{(i)}$ decrypts to β when the above algorithm is run with the function $f^{(i)}$.

Unfortunately, the communication channels are very noisy, so each of the n bits in each of the k transmissions is *independently* corrupted (i.e., flipped to its complement) with probability $1/4$. Thus no single receiver on his or her own is likely to be able to decrypt the message correctly. Show, however, that if k is large enough as a function of n , then the k receivers can jointly reconstruct the plain-text message in the following way. They get together, and without revealing any of the $\alpha^{(i)}$ or the $f^{(i)}$, they interactively run an algorithm that will produce the correct β with probability at least $9/10$. (How large do you need k to be in your algorithm?)

17. Consider the following simple model of gambling in the presence of bad odds. At the beginning, your net profit is 0. You play for a sequence of n rounds; and in each round, your net profit increases by 1 with probability $1/3$, and decreases by 1 with probability $2/3$.

Show that the expected number of steps in which your net profit is positive can be upper-bounded by an absolute constant, independent of the value of n .

18. In this problem, we will consider the following simple randomized algorithm for the Vertex Cover Algorithm.

```

Start with  $S = \emptyset$ 
While  $S$  is not a vertex cover,
  Select an edge  $e$  not covered by  $S$ 
  Select one end of  $e$  at random (each end equally likely)
  Add the selected node to  $S$ 
Endwhile

```

We will be interested in the expected cost of a vertex cover selected by this algorithm.

- (a) Is this algorithm a c -approximation algorithm for the Minimum Weight Vertex Cover Problem for some constant c ? Prove your answer.
- (b) Is this algorithm a c -approximation algorithm for the Minimum Cardinality Vertex Cover Problem for some constant c ? Prove your answer.

(Hint: For an edge, let p_e denote the probability that edge e is selected as an uncovered edge in this algorithm. Can you express the expected value of the solution in terms of these probabilities? To bound the value of an optimal solution in terms of the p_e probabilities, try to bound the sum of the probabilities for the edges incident to a given vertex v , namely, $\sum_{e \text{ incident to } v} p_e$.)

Notes and Further Reading

The use of randomization in algorithms is an active research area; the books by Motwani and Raghavan (1995) and Mitzenmacher and Upfal (2005) are devoted to this topic. As the contents of this chapter make clear, the types of probabilistic arguments used in the study of basic randomized algorithms often have a discrete, combinatorial flavor; one can get background in this style of probabilistic analysis from the book by Feller (1957).

The use of randomization for contention resolution is common in many systems and networking applications. Ethernet-style shared communication media, for example, use randomized *backoff* protocols to reduce the number of collisions among different senders; see the book by Bertsekas and Gallager (1992) for a discussion of this topic.

The randomized algorithm for the Minimum-Cut Problem described in the text is due to Karger, and after further optimizations due to Karger and Stein (1996), it has become one of the most efficient approaches to the minimum cut problem. A number of further extensions and applications of the algorithm appear in Karger's (1995) Ph.D. thesis.

The approximation algorithm for MAX 3-SAT is due to Johnson (1974), in a paper that contains a number of early approximation algorithms for NP-hard problems. The surprising punch line to that section—that every instance of 3-SAT has an assignment satisfying at least 7/8 of the clauses—is an example of the *probabilistic method*, whereby a combinatorial structure with a desired property is shown to exist simply by arguing that a random structure has the property with positive probability. This has grown into a highly refined