

# Lecture 7

- Dynamic programming.

Dynamic programming is a design strategy that is useful when the underlying problem demonstrates a kind of optimal substructure, but the sub-problems are *overlapping*. To see what we mean by this, consider merge sort, which deploys divide-n-conquer. The problem of sorting an array  $A[1, \dots, n]$  is split in to the two sub-problems of sorting  $A[1, \dots, n/2]$  and  $[n/2 + 1, \dots, n]$ . These two sub-problems do not overlap — the sub-array  $A[1, \dots, n/2]$ , which characterizes one sub-problem is completely distinct from the sub-array  $A[n/2 + 1, \dots, n]$ , which characterizes the other sub-problem.

Now consider the problem of computing a Fibonacci number, which is expressed by the following recurrence for  $f: \mathbb{Z}_0^+ \rightarrow \mathbb{Z}_0^+$ .

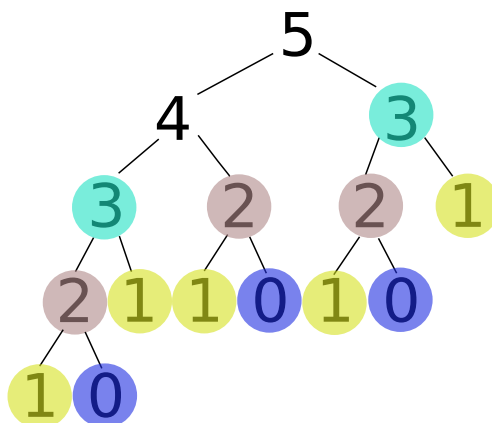
$$f: x \mapsto \begin{cases} 0 & \text{if } x = 0 \\ 1 & \text{if } x = 1 \\ f(x-1) + f(x-2) & \text{otherwise} \end{cases}$$

Following is a straightforward encoding of the above recurrence using a recursive algorithm.

```
NAIVEFIBO( $x$ )  
  1 if  $x = 0$  then return 0  
  2 if  $x = 1$  then return 1  
  3 return NAIVEFIBO( $x - 1$ ) + NAIVEFIBO( $x - 2$ )
```

Now suppose we draw a recursion tree to see all the calls that are made

when we invoke `NAIVEFIBO` with some argument, say  $x = 5$ . The following picture shows all the calls that result. The overlapping sub-problems are shown in circles of the same colour. For example, we have three distinct calls to `NAIVEFIBO` with argument 2.



Adopting the total number of calls to `NAIVEFIBO` as a measure of time-efficiency, we can prove that the total number of calls is  $\Omega(2^x)$  on input  $x$ .

One way to avoid such redundant “overlapping” sub-problems is to *memoize* intermediate results. The word “memoize” is as in writing a memo to yourself for later use. Following is a version of an algorithm that uses memoization.

<pre> MEMOFIBO(<math>x</math>) 1 Allocate array <math>F[0, \dots, x]</math> 2 <math>F[0] \leftarrow 0, F[1] \leftarrow 1</math> 3 <b>foreach</b> <math>i</math> from 2 to <math>x</math> <b>do</b> 4   <math>F[i] \leftarrow \text{NIL}</math> 5 <math>\text{COMPUTEFIBO}(F, x)</math> 6 <b>return</b> <math>F[x]</math> </pre>	<pre> COMPUTEFIBO(<math>F, x</math>) 11 <b>if</b> <math>F[x] = \text{NIL}</math> <b>then</b> 12   <math>F[x] \leftarrow \text{COMPUTEFIBO}(F, x - 1) +</math>       <math>\text{COMPUTEFIBO}(F, x - 2)</math> 13 <b>return</b> <math>F[x]</math> </pre>
---	---

The above version guarantees that each sub-problem is solved once only. Thus, the total number of calls to `COMPUTEFIBO` that result from a call `MEMOFIBO( $x$ )` is  $\Theta(x)$ .

A final improvement we could do is to not use recursion at all, and simply use iteration, but go “bottom-up” instead of top-down. This is exactly dynamic

programming. That is, dynamic programming is memoization to address overlapping sub-problems, and doing the computations bottom-up instead of top-down.

```

DPFIBO( $x$ )
1 Allocate array  $F[0, \dots, x]$ 
2  $F[0] \leftarrow 0, F[1] \leftarrow 1$ 
3 foreach  $i$  from 2 to  $x$  do
4      $F[i] \leftarrow F[i-1] + F[i-2]$ 
5 return  $F[x]$ 

```

The asymptotic running-time remains unchanged from the prior version that is top-down with memoization. But we do save some overhead in comparison to that version from not having to make any recursive calls.

**Optimal substructure** The specification of a Fibonacci number we adopt above is as a recurrence. This actually eases the first step for us, which is of recognizing that the problem possesses optimal substructure. This is a key first step, and indeed, a valuable skill in the design of algorithms — recognition that a problem possesses some kind of optimal substructure. The way to do this is to go “top-down” — ask whether, given solutions to sub-problems, we are able to compose them somehow to yield a solution to the bigger problem.

Consider, as an example, interval scheduling. We now reproduce the discussions from CLRS on its optimal substructure, and the manner in which that can be exploited to get an efficient algorithm based on dynamic programming. CLRS calls the problem “activity selection.”

After that, we include from CLRS some of their additional material on dynamic programming, specifically, the problems of rod selection and matrix multiplication. Finally, we include the ingenious Floyd-Warshall’s algorithm for all pairs shortest paths.

## 16.1 An activity-selection problem

Our first example is the problem of scheduling several competing activities that require exclusive use of a common resource, with a goal of selecting a maximum-size set of mutually compatible activities. Suppose we have a set  $S = \{a_1, a_2, \dots, a_n\}$  of  $n$  proposed **activities** that wish to use a resource, such as a lecture hall, which can be used by only one activity at a time. Each activity  $a_i$  has a **start time**  $s_i$  and a **finish time**  $f_i$ , where  $0 \leq s_i < f_i < \infty$ . If selected, activity  $a_i$  takes place during the half-open time interval  $[s_i, f_i)$ . Activities  $a_i$  and  $a_j$  are **compatible** if the intervals  $[s_i, f_i)$  and  $[s_j, f_j)$  do not overlap (i.e.,  $a_i$  and  $a_j$  are compatible if  $s_i \geq f_j$  or  $s_j \geq f_i$ ). The **activity-selection problem** is to select a maximum-size subset of mutually compatible activities. For example, consider the following set  $S$  of activities, which we have sorted in monotonically increasing order of finish time:

$i$	1	2	3	4	5	6	7	8	9	10	11
$s_i$	1	3	0	5	3	5	6	8	8	2	12
$f_i$	4	5	6	7	8	9	10	11	12	13	14

(We shall see shortly why it is advantageous to consider activities in sorted order.) For this example, the subset  $\{a_3, a_9, a_{11}\}$  consists of mutually compatible activities. It is not a maximal subset, however, since the subset  $\{a_1, a_4, a_8, a_{11}\}$  is larger. In fact,  $\{a_1, a_4, a_8, a_{11}\}$  is a largest subset of mutually compatible activities; another largest subset is  $\{a_2, a_4, a_9, a_{11}\}$ .

We shall solve this problem in several steps. We start by formulating a dynamic-programming solution to this problem in which we combine optimal solutions to two subproblems to form an optimal solution to the original problem. We consider several choices when determining which subproblems to use in an optimal solution. We shall then observe that we need only consider one choice—the greedy choice—and that when we make the greedy choice, one of the subproblems is guaranteed to be empty, so that only one nonempty subproblem remains. Based on these observations, we shall develop a recursive greedy algorithm to solve the activity-scheduling problem. We shall complete the process of developing a greedy solution by converting the recursive algorithm to an iterative one. Although the steps we shall go through in this section are more involved than is typical for the development of a greedy algorithm, they illustrate the relationship of greedy algorithms and dynamic programming.

### The optimal substructure of the activity-selection problem

As mentioned above, we start by developing a dynamic-programming solution to the activity-selection problem. As in Chapter 15, our first step is to find the optimal

substructure and then use it to construct an optimal solution to the problem from optimal solutions to subproblems.

We saw in Chapter 15 that we need to define an appropriate space of subproblems. Let us start by defining sets

$$S_{ij} = \{a_k \in S : f_i \leq s_k < f_k \leq s_j\} ,$$

so that  $S_{ij}$  is the subset of activities in  $S$  that can start after activity  $a_i$  finishes and finish before activity  $a_j$  starts. In fact,  $S_{ij}$  consists of all activities that are compatible with  $a_i$  and  $a_j$  and are also compatible with all activities that finish no later than  $a_i$  finishes and all activities that start no earlier than  $a_j$  starts. In order to represent the entire problem, we add fictitious activities  $a_0$  and  $a_{n+1}$  and adopt the conventions that  $f_0 = 0$  and  $s_{n+1} = \infty$ . Then  $S = S_{0,n+1}$ , and the ranges for  $i$  and  $j$  are given by  $0 \leq i, j \leq n + 1$ .

We can further restrict the ranges of  $i$  and  $j$  as follows. Let us assume that the activities are sorted in monotonically increasing order of finish time:

$$f_0 \leq f_1 \leq f_2 \leq \cdots \leq f_n < f_{n+1} . \quad (16.1)$$

We claim that  $S_{ij} = \emptyset$  whenever  $i \geq j$ . Why? Suppose that there exists an activity  $a_k \in S_{ij}$  for some  $i \geq j$ , so that  $a_i$  follows  $a_j$  in the sorted order. Then we would have  $f_i \leq s_k < f_k \leq s_j < f_j$ . Thus,  $f_i < f_j$ , which contradicts our assumption that  $a_i$  follows  $a_j$  in the sorted order. We can conclude that, assuming that we have sorted the activities in monotonically increasing order of finish time, our space of subproblems is to select a maximum-size subset of mutually compatible activities from  $S_{ij}$ , for  $0 \leq i < j \leq n + 1$ , knowing that all other  $S_{ij}$  are empty.

To see the substructure of the activity-selection problem, consider some non-empty subproblem  $S_{ij}$ ,<sup>1</sup> and suppose that a solution to  $S_{ij}$  includes some activity  $a_k$ , so that  $f_i \leq s_k < f_k \leq s_j$ . Using activity  $a_k$  generates two subproblems,  $S_{ik}$  (activities that start after  $a_i$  finishes and finish before  $a_k$  starts) and  $S_{kj}$  (activities that start after  $a_k$  finishes and finish before  $a_j$  starts), each of which consists of a subset of the activities in  $S_{ij}$ . Our solution to  $S_{ij}$  is the union of the solutions to  $S_{ik}$  and  $S_{kj}$ , along with the activity  $a_k$ . Thus, the number of activities in our solution to  $S_{ij}$  is the size of our solution to  $S_{ik}$ , plus the size of our solution to  $S_{kj}$ , plus one (for  $a_k$ ).

The optimal substructure of this problem is as follows. Suppose now that an optimal solution  $A_{ij}$  to  $S_{ij}$  includes activity  $a_k$ . Then the solutions  $A_{ik}$  to  $S_{ik}$  and  $A_{kj}$  to  $S_{kj}$  used within this optimal solution to  $S_{ij}$  must be optimal as well. The usual cut-and-paste argument applies. If we had a solution  $A'_{ik}$  to  $S_{ik}$  that included

---

<sup>1</sup>We will sometimes speak of the sets  $S_{ij}$  as subproblems rather than just sets of activities. It will always be clear from the context whether we are referring to  $S_{ij}$  as a set of activities or the subproblem whose input is that set.

more activities than  $A_{ik}$ , we could cut out  $A_{ik}$  from  $A_{ij}$  and paste in  $A'_{ik}$ , thus producing another solution to  $S_{ij}$  with more activities than  $A_{ij}$ . Because we assumed that  $A_{ij}$  is an optimal solution, we have derived a contradiction. Similarly, if we had a solution  $A'_{kj}$  to  $S_{kj}$  with more activities than  $A_{kj}$ , we could replace  $A_{kj}$  by  $A'_{kj}$  to produce a solution to  $S_{ij}$  with more activities than  $A_{ij}$ .

Now we use our optimal substructure to show that we can construct an optimal solution to the problem from optimal solutions to subproblems. We have seen that any solution to a nonempty subproblem  $S_{ij}$  includes some activity  $a_k$ , and that any optimal solution contains within it optimal solutions to subproblem instances  $S_{ik}$  and  $S_{kj}$ . Thus, we can build a maximum-size subset of mutually compatible activities in  $S_{ij}$  by splitting the problem into two subproblems (finding maximum-size subsets of mutually compatible activities in  $S_{ik}$  and  $S_{kj}$ ), finding maximum-size subsets  $A_{ik}$  and  $A_{kj}$  of mutually compatible activities for these subproblems, and forming our maximum-size subset  $A_{ij}$  of mutually compatible activities as

$$A_{ij} = A_{ik} \cup \{a_k\} \cup A_{kj} . \quad (16.2)$$

An optimal solution to the entire problem is a solution to  $S_{0,n+1}$ .

### A recursive solution

The second step in developing a dynamic-programming solution is to recursively define the value of an optimal solution. For the activity-selection problem, we let  $c[i, j]$  be the number of activities in a maximum-size subset of mutually compatible activities in  $S_{ij}$ . We have  $c[i, j] = 0$  whenever  $S_{ij} = \emptyset$ ; in particular,  $c[i, j] = 0$  for  $i \geq j$ .

Now consider a nonempty subset  $S_{ij}$ . As we have seen, if  $a_k$  is used in a maximum-size subset of mutually compatible activities of  $S_{ij}$ , we also use maximum-size subsets of mutually compatible activities for the subproblems  $S_{ik}$  and  $S_{kj}$ . Using equation (16.2), we have the recurrence

$$c[i, j] = c[i, k] + c[k, j] + 1 .$$

This recursive equation assumes that we know the value of  $k$ , which we do not. There are  $j - i - 1$  possible values for  $k$ , namely  $k = i + 1, \dots, j - 1$ . Since the maximum-size subset of  $S_{ij}$  must use one of these values for  $k$ , we check them all to find the best. Thus, our full recursive definition of  $c[i, j]$  becomes

$$c[i, j] = \begin{cases} 0 & \text{if } S_{ij} = \emptyset , \\ \max_{\substack{i < k < j \\ a_k \in S_{ij}}} \{c[i, k] + c[k, j] + 1\} & \text{if } S_{ij} \neq \emptyset . \end{cases} \quad (16.3)$$

---

## 15      Dynamic Programming

Dynamic programming, like the divide-and-conquer method, solves problems by combining the solutions to subproblems. (“Programming” in this context refers to a tabular method, not to writing computer code.) As we saw in Chapters 2 and 4, divide-and-conquer algorithms partition the problem into disjoint subproblems, solve the subproblems recursively, and then combine their solutions to solve the original problem. In contrast, dynamic programming applies when the subproblems overlap—that is, when subproblems share subsubproblems. In this context, a divide-and-conquer algorithm does more work than necessary, repeatedly solving the common subsubproblems. A dynamic-programming algorithm solves each subsubproblem just once and then saves its answer in a table, thereby avoiding the work of recomputing the answer every time it solves each subsubproblem.

We typically apply dynamic programming to *optimization problems*. Such problems can have many possible solutions. Each solution has a value, and we wish to find a solution with the optimal (minimum or maximum) value. We call such a solution *an* optimal solution to the problem, as opposed to *the* optimal solution, since there may be several solutions that achieve the optimal value.

When developing a dynamic-programming algorithm, we follow a sequence of four steps:

1. Characterize the structure of an optimal solution.
2. Recursively define the value of an optimal solution.
3. Compute the value of an optimal solution, typically in a bottom-up fashion.
4. Construct an optimal solution from computed information.

Steps 1–3 form the basis of a dynamic-programming solution to a problem. If we need only the value of an optimal solution, and not the solution itself, then we can omit step 4. When we do perform step 4, we sometimes maintain additional information during step 3 so that we can easily construct an optimal solution.

The sections that follow use the dynamic-programming method to solve some optimization problems. Section 15.1 examines the problem of cutting a rod into

rods of smaller length in a way that maximizes their total value. Section 15.2 asks how we can multiply a chain of matrices while performing the fewest total scalar multiplications. Given these examples of dynamic programming, Section 15.3 discusses two key characteristics that a problem must have for dynamic programming to be a viable solution technique. Section 15.4 then shows how to find the longest common subsequence of two sequences via dynamic programming. Finally, Section 15.5 uses dynamic programming to construct binary search trees that are optimal, given a known distribution of keys to be looked up.

15.1 Rod cutting

Our first example uses dynamic programming to solve a simple problem in deciding where to cut steel rods. Serling Enterprises buys long steel rods and cuts them into shorter rods, which it then sells. Each cut is free. The management of Serling Enterprises wants to know the best way to cut up the rods.

We assume that we know, for  $i = 1, 2, \dots$ , the price  $p_i$  in dollars that Serling Enterprises charges for a rod of length  $i$  inches. Rod lengths are always an integral number of inches. Figure 15.1 gives a sample price table.

The **rod-cutting problem** is the following. Given a rod of length  $n$  inches and a table of prices  $p_i$  for  $i = 1, 2, \dots, n$ , determine the maximum revenue  $r_n$  obtainable by cutting up the rod and selling the pieces. Note that if the price  $p_n$  for a rod of length  $n$  is large enough, an optimal solution may require no cutting at all.

Consider the case when  $n = 4$ . Figure 15.2 shows all the ways to cut up a rod of 4 inches in length, including the way with no cuts at all. We see that cutting a 4-inch rod into two 2-inch pieces produces revenue  $p_2 + p_2 = 5 + 5 = 10$ , which is optimal.

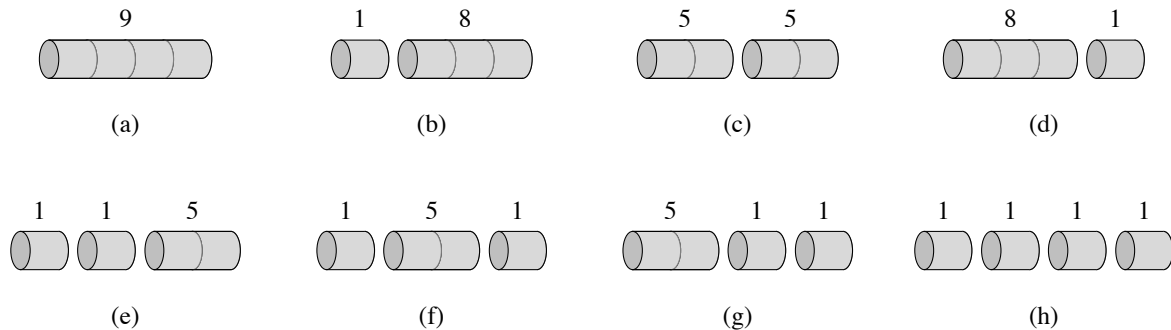
We can cut up a rod of length  $n$  in  $2^{n-1}$  different ways, since we have an independent option of cutting, or not cutting, at distance  $i$  inches from the left end,

length $i$	1	2	3	4	5	6	7	8	9	10
price $p_i$	1	5	8	9	10	17	17	20	24	30

**Figure 15.1** A sample price table for rods. Each rod of length  $i$  inches earns the company  $p_i$  dollars of revenue.

Copyright © 2009, MIT Press. All rights reserved.





**Figure 15.2** The 8 possible ways of cutting up a rod of length 4. Above each piece is the value of that piece, according to the sample price chart of Figure 15.1. The optimal strategy is part (c)—cutting the rod into two pieces of length 2—which has total value 10.

for  $i = 1, 2, \dots, n - 1$ .<sup>1</sup> We denote a decomposition into pieces using ordinary additive notation, so that  $7 = 2 + 2 + 3$  indicates that a rod of length 7 is cut into three pieces—two of length 2 and one of length 3. If an optimal solution cuts the rod into  $k$  pieces, for some  $1 \leq k \leq n$ , then an optimal decomposition

$$n = i_1 + i_2 + \dots + i_k$$

of the rod into pieces of lengths  $i_1, i_2, \dots, i_k$  provides maximum corresponding revenue

$$r_n = p_{i_1} + p_{i_2} + \dots + p_{i_k}.$$

For our sample problem, we can determine the optimal revenue figures  $r_i$ , for  $i = 1, 2, \dots, 10$ , by inspection, with the corresponding optimal decompositions

<sup>1</sup>If we required the pieces to be cut in order of nondecreasing size, there would be fewer ways to consider. For  $n = 4$ , we would consider only 5 such ways: parts (a), (b), (c), (e), and (h) in Figure 15.2. The number of ways is called the *partition function*; it is approximately equal to  $e^{\pi\sqrt{2n/3}}/4n\sqrt{3}$ . This quantity is less than  $2^{n-1}$ , but still much greater than any polynomial in  $n$ . We shall not pursue this line of inquiry further, however.

$$\begin{aligned}
r_1 &= 1 && \text{from solution } 1 = 1 \quad (\text{no cuts}) , \\
r_2 &= 5 && \text{from solution } 2 = 2 \quad (\text{no cuts}) , \\
r_3 &= 8 && \text{from solution } 3 = 3 \quad (\text{no cuts}) , \\
r_4 &= 10 && \text{from solution } 4 = 2 + 2 , \\
r_5 &= 13 && \text{from solution } 5 = 2 + 3 , \\
r_6 &= 17 && \text{from solution } 6 = 6 \quad (\text{no cuts}) , \\
r_7 &= 18 && \text{from solution } 7 = 1 + 6 \text{ or } 7 = 2 + 2 + 3 , \\
r_8 &= 22 && \text{from solution } 8 = 2 + 6 , \\
r_9 &= 25 && \text{from solution } 9 = 3 + 6 , \\
r_{10} &= 30 && \text{from solution } 10 = 10 \quad (\text{no cuts}) .
\end{aligned}$$

More generally, we can frame the values  $r_n$  for  $n \geq 1$  in terms of optimal revenues from shorter rods:

$$r_n = \max(p_n, r_1 + r_{n-1}, r_2 + r_{n-2}, \dots, r_{n-1} + r_1) . \quad (15.1)$$

The first argument,  $p_n$ , corresponds to making no cuts at all and selling the rod of length  $n$  as is. The other  $n - 1$  arguments to max correspond to the maximum revenue obtained by making an initial cut of the rod into two pieces of size  $i$  and  $n - i$ , for each  $i = 1, 2, \dots, n - 1$ , and then optimally cutting up those pieces further, obtaining revenues  $r_i$  and  $r_{n-i}$  from those two pieces. Since we don't know ahead of time which value of  $i$  optimizes revenue, we have to consider all possible values for  $i$  and pick the one that maximizes revenue. We also have the option of picking no  $i$  at all if we can obtain more revenue by selling the rod uncut.

Note that to solve the original problem of size  $n$ , we solve smaller problems of the same type, but of smaller sizes. Once we make the first cut, we may consider the two pieces as independent instances of the rod-cutting problem. The overall optimal solution incorporates optimal solutions to the two related subproblems, maximizing revenue from each of those two pieces. We say that the rod-cutting problem exhibits **optimal substructure**: optimal solutions to a problem incorporate optimal solutions to related subproblems, which we may solve independently.

In a related, but slightly simpler, way to arrange a recursive structure for the rod-cutting problem, we view a decomposition as consisting of a first piece of length  $i$  cut off the left-hand end, and then a right-hand remainder of length  $n - i$ . Only the remainder, and not the first piece, may be further divided. We may view every decomposition of a length- $n$  rod in this way: as a first piece followed by some decomposition of the remainder. When doing so, we can couch the solution with no cuts at all as saying that the first piece has size  $i = n$  and revenue  $p_n$  and that the remainder has size 0 with corresponding revenue  $r_0 = 0$ . We thus obtain the following simpler version of equation (15.1):

$$r_n = \max_{1 \leq i \leq n} (p_i + r_{n-i}) . \quad (15.2)$$

In this formulation, an optimal solution embodies the solution to only *one* related subproblem—the remainder—rather than two.

### Recursive top-down implementation

The following procedure implements the computation implicit in equation (15.2) in a straightforward, top-down, recursive manner.

CUT-ROD( $p, n$ )

```

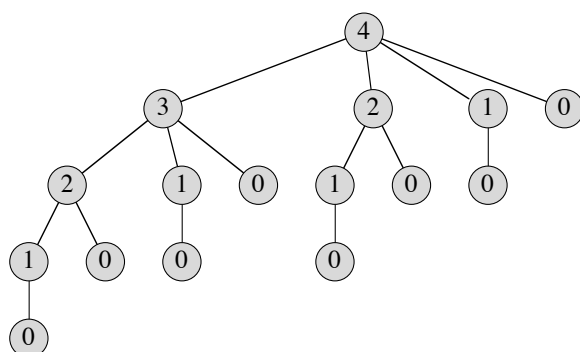
1  if  $n == 0$ 
2      return 0
3   $q = -\infty$ 
4  for  $i = 1$  to  $n$ 
5       $q = \max(q, p[i] + \text{CUT-ROD}(p, n - i))$ 
6  return  $q$ 
```

Procedure CUT-ROD takes as input an array  $p[1..n]$  of prices and an integer  $n$ , and it returns the maximum revenue possible for a rod of length  $n$ . If  $n = 0$ , no revenue is possible, and so CUT-ROD returns 0 in line 2. Line 3 initializes the maximum revenue  $q$  to  $-\infty$ , so that the **for** loop in lines 4–5 correctly computes  $q = \max_{1 \leq i \leq n} (p_i + \text{CUT-ROD}(p, n - i))$ ; line 6 then returns this value. A simple induction on  $n$  proves that this answer is equal to the desired answer  $r_n$ , using equation (15.2).

If you were to code up CUT-ROD in your favorite programming language and run it on your computer, you would find that once the input size becomes moderately large, your program would take a long time to run. For  $n = 40$ , you would find that your program takes at least several minutes, and most likely more than an hour. In fact, you would find that each time you increase  $n$  by 1, your program's running time would approximately double.

Why is CUT-ROD so inefficient? The problem is that CUT-ROD calls itself recursively over and over again with the same parameter values; it solves the same subproblems repeatedly. Figure 15.3 illustrates what happens for  $n = 4$ : CUT-ROD( $p, n$ ) calls CUT-ROD( $p, n - i$ ) for  $i = 1, 2, \dots, n$ . Equivalently, CUT-ROD( $p, n$ ) calls CUT-ROD( $p, j$ ) for each  $j = 0, 1, \dots, n - 1$ . When this process unfolds recursively, the amount of work done, as a function of  $n$ , grows explosively.

To analyze the running time of CUT-ROD, let  $T(n)$  denote the total number of calls made to CUT-ROD when called with its second parameter equal to  $n$ . This expression equals the number of nodes in a subtree whose root is labeled  $n$  in the recursion tree. The count includes the initial call at its root. Thus,  $T(0) = 1$  and



**Figure 15.3** The recursion tree showing recursive calls resulting from a call  $\text{CUT-ROD}(p, n)$  for  $n = 4$ . Each node label gives the size  $n$  of the corresponding subproblem, so that an edge from a parent with label  $s$  to a child with label  $t$  corresponds to cutting off an initial piece of size  $s - t$  and leaving a remaining subproblem of size  $t$ . A path from the root to a leaf corresponds to one of the  $2^{n-1}$  ways of cutting up a rod of length  $n$ . In general, this recursion tree has  $2^n$  nodes and  $2^{n-1}$  leaves.

$$T(n) = 1 + \sum_{j=0}^{n-1} T(j). \quad (15.3)$$

The initial 1 is for the call at the root, and the term  $T(j)$  counts the number of calls (including recursive calls) due to the call  $\text{CUT-ROD}(p, n - i)$ , where  $j = n - i$ . As Exercise 15.1-1 asks you to show,

$$T(n) = 2^n, \quad (15.4)$$

and so the running time of  $\text{CUT-ROD}$  is exponential in  $n$ .

In retrospect, this exponential running time is not so surprising.  $\text{CUT-ROD}$  explicitly considers all the  $2^{n-1}$  possible ways of cutting up a rod of length  $n$ . The tree of recursive calls has  $2^{n-1}$  leaves, one for each possible way of cutting up the rod. The labels on the simple path from the root to a leaf give the sizes of each remaining right-hand piece before making each cut. That is, the labels give the corresponding cut points, measured from the right-hand end of the rod.

### Using dynamic programming for optimal rod cutting

We now show how to convert  $\text{CUT-ROD}$  into an efficient algorithm, using dynamic programming.

The dynamic-programming method works as follows. Having observed that a naive recursive solution is inefficient because it solves the same subproblems repeatedly, we arrange for each subproblem to be solved only *once*, saving its solution. If we need to refer to this subproblem's solution again later, we can just look it

up, rather than recompute it. Dynamic programming thus uses additional memory to save computation time; it serves an example of a *time-memory trade-off*. The savings may be dramatic: an exponential-time solution may be transformed into a polynomial-time solution. A dynamic-programming approach runs in polynomial time when the number of *distinct* subproblems involved is polynomial in the input size and we can solve each such subproblem in polynomial time.

There are usually two equivalent ways to implement a dynamic-programming approach. We shall illustrate both of them with our rod-cutting example.

The first approach is *top-down with memoization*.<sup>2</sup> In this approach, we write the procedure recursively in a natural manner, but modified to save the result of each subproblem (usually in an array or hash table). The procedure now first checks to see whether it has previously solved this subproblem. If so, it returns the saved value, saving further computation at this level; if not, the procedure computes the value in the usual manner. We say that the recursive procedure has been *memoized*; it “remembers” what results it has computed previously.

The second approach is the *bottom-up method*. This approach typically depends on some natural notion of the “size” of a subproblem, such that solving any particular subproblem depends only on solving “smaller” subproblems. We sort the subproblems by size and solve them in size order, smallest first. When solving a particular subproblem, we have already solved all of the smaller subproblems its solution depends upon, and we have saved their solutions. We solve each subproblem only once, and when we first see it, we have already solved all of its prerequisite subproblems.

These two approaches yield algorithms with the same asymptotic running time, except in unusual circumstances where the top-down approach does not actually recurse to examine all possible subproblems. The bottom-up approach often has much better constant factors, since it has less overhead for procedure calls.

Here is the pseudocode for the top-down CUT-ROD procedure, with memoization added:

MEMOIZED-CUT-ROD( $p, n$ )

```

1  let  $r[0 \dots n]$  be a new array
2  for  $i = 0$  to  $n$ 
3       $r[i] = -\infty$ 
4  return MEMOIZED-CUT-ROD-AUX( $p, n, r$ )
```

---

<sup>2</sup>This is not a misspelling. The word really is *memoization*, not *memorization*. *Memoization* comes from *memo*, since the technique consists of recording a value so that we can look it up later.

MEMOIZED-CUT-ROD-AUX( $p, n, r$ )

```

1  if  $r[n] \geq 0$ 
2      return  $r[n]$ 
3  if  $n == 0$ 
4       $q = 0$ 
5  else  $q = -\infty$ 
6      for  $i = 1$  to  $n$ 
7           $q = \max(q, p[i] + \text{MEMOIZED-CUT-ROD-AUX}(p, n - i, r))$ 
8   $r[n] = q$ 
9  return  $q$ 

```

Here, the main procedure MEMOIZED-CUT-ROD initializes a new auxiliary array  $r[0..n]$  with the value  $-\infty$ , a convenient choice with which to denote “unknown.” (Known revenue values are always nonnegative.) It then calls its helper routine, MEMOIZED-CUT-ROD-AUX.

The procedure MEMOIZED-CUT-ROD-AUX is just the memoized version of our previous procedure, CUT-ROD. It first checks in line 1 to see whether the desired value is already known and, if it is, then line 2 returns it. Otherwise, lines 3–7 compute the desired value  $q$  in the usual manner, line 8 saves it in  $r[n]$ , and line 9 returns it.

The bottom-up version is even simpler:

BOTTOM-UP-CUT-ROD( $p, n$ )

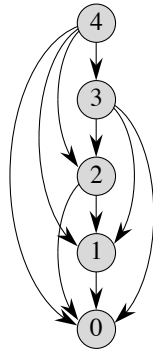
```

1  let  $r[0..n]$  be a new array
2   $r[0] = 0$ 
3  for  $j = 1$  to  $n$ 
4       $q = -\infty$ 
5      for  $i = 1$  to  $j$ 
6           $q = \max(q, p[i] + r[j - i])$ 
7       $r[j] = q$ 
8  return  $r[n]$ 

```

For the bottom-up dynamic-programming approach, BOTTOM-UP-CUT-ROD uses the natural ordering of the subproblems: a subproblem of size  $i$  is “smaller” than a subproblem of size  $j$  if  $i < j$ . Thus, the procedure solves subproblems of sizes  $j = 0, 1, \dots, n$ , in that order.

Line 1 of procedure BOTTOM-UP-CUT-ROD creates a new array  $r[0..n]$  in which to save the results of the subproblems, and line 2 initializes  $r[0]$  to 0, since a rod of length 0 earns no revenue. Lines 3–6 solve each subproblem of size  $j$ , for  $j = 1, 2, \dots, n$ , in order of increasing size. The approach used to solve a problem of a particular size  $j$  is the same as that used by CUT-ROD, except that line 6 now



**Figure 15.4** The subproblem graph for the rod-cutting problem with  $n = 4$ . The vertex labels give the sizes of the corresponding subproblems. A directed edge  $(x, y)$  indicates that we need a solution to subproblem  $y$  when solving subproblem  $x$ . This graph is a reduced version of the tree of Figure 15.3, in which all nodes with the same label are collapsed into a single vertex and all edges go from parent to child.

directly references array entry  $r[j - i]$  instead of making a recursive call to solve the subproblem of size  $j - i$ . Line 7 saves in  $r[j]$  the solution to the subproblem of size  $j$ . Finally, line 8 returns  $r[n]$ , which equals the optimal value  $r_n$ .

The bottom-up and top-down versions have the same asymptotic running time. The running time of procedure BOTTOM-UP-CUT-ROD is  $\Theta(n^2)$ , due to its doubly-nested loop structure. The number of iterations of its inner **for** loop, in lines 5–6, forms an arithmetic series. The running time of its top-down counterpart, MEMOIZED-CUT-ROD, is also  $\Theta(n^2)$ , although this running time may be a little harder to see. Because a recursive call to solve a previously solved subproblem returns immediately, MEMOIZED-CUT-ROD solves each subproblem just once. It solves subproblems for sizes  $0, 1, \dots, n$ . To solve a subproblem of size  $n$ , the **for** loop of lines 6–7 iterates  $n$  times. Thus, the total number of iterations of this **for** loop, over all recursive calls of MEMOIZED-CUT-ROD, forms an arithmetic series, giving a total of  $\Theta(n^2)$  iterations, just like the inner **for** loop of BOTTOM-UP-CUT-ROD. (We actually are using a form of aggregate analysis here. We shall see aggregate analysis in detail in Section 17.1.)

### Subproblem graphs

When we think about a dynamic-programming problem, we should understand the set of subproblems involved and how subproblems depend on one another.

The **subproblem graph** for the problem embodies exactly this information. Figure 15.4 shows the subproblem graph for the rod-cutting problem with  $n = 4$ . It is a directed graph, containing one vertex for each distinct subproblem. The sub-

problem graph has a directed edge from the vertex for subproblem  $x$  to the vertex for subproblem  $y$  if determining an optimal solution for subproblem  $x$  involves directly considering an optimal solution for subproblem  $y$ . For example, the subproblem graph contains an edge from  $x$  to  $y$  if a top-down recursive procedure for solving  $x$  directly calls itself to solve  $y$ . We can think of the subproblem graph as a “reduced” or “collapsed” version of the recursion tree for the top-down recursive method, in which we coalesce all nodes for the same subproblem into a single vertex and direct all edges from parent to child.

The bottom-up method for dynamic programming considers the vertices of the subproblem graph in such an order that we solve the subproblems  $y$  adjacent to a given subproblem  $x$  before we solve subproblem  $x$ . (Recall from Section B.4 that the adjacency relation is not necessarily symmetric.) Using the terminology from Chapter 22, in a bottom-up dynamic-programming algorithm, we consider the vertices of the subproblem graph in an order that is a “reverse topological sort,” or a “topological sort of the transpose” (see Section 22.4) of the subproblem graph. In other words, no subproblem is considered until all of the subproblems it depends upon have been solved. Similarly, using notions from the same chapter, we can view the top-down method (with memoization) for dynamic programming as a “depth-first search” of the subproblem graph (see Section 22.3).

The size of the subproblem graph  $G = (V, E)$  can help us determine the running time of the dynamic programming algorithm. Since we solve each subproblem just once, the running time is the sum of the times needed to solve each subproblem. Typically, the time to compute the solution to a subproblem is proportional to the degree (number of outgoing edges) of the corresponding vertex in the subproblem graph, and the number of subproblems is equal to the number of vertices in the subproblem graph. In this common case, the running time of dynamic programming is linear in the number of vertices and edges.

### Reconstructing a solution

Our dynamic-programming solutions to the rod-cutting problem return the value of an optimal solution, but they do not return an actual solution: a list of piece sizes. We can extend the dynamic-programming approach to record not only the optimal *value* computed for each subproblem, but also a *choice* that led to the optimal value. With this information, we can readily print an optimal solution.

Here is an extended version of BOTTOM-UP-CUT-ROD that computes, for each rod size  $j$ , not only the maximum revenue  $r_j$ , but also  $s_j$ , the optimal size of the first piece to cut off:



EXTENDED-BOTTOM-UP-CUT-ROD( $p, n$ )

```

1  let  $r[0..n]$  and  $s[1..n]$  be new arrays
2   $r[0] = 0$ 
3  for  $j = 1$  to  $n$ 
4       $q = -\infty$ 
5      for  $i = 1$  to  $j$ 
6          if  $q < p[i] + r[j - i]$ 
7               $q = p[i] + r[j - i]$ 
8               $s[j] = i$ 
9   $r[j] = q$ 
10 return  $r$  and  $s$ 
```

This procedure is similar to BOTTOM-UP-CUT-ROD, except that it creates the array  $s$  in line 1, and it updates  $s[j]$  in line 8 to hold the optimal size  $i$  of the first piece to cut off when solving a subproblem of size  $j$ .

The following procedure takes a price table  $p$  and a rod size  $n$ , and it calls EXTENDED-BOTTOM-UP-CUT-ROD to compute the array  $s[1..n]$  of optimal first-piece sizes and then prints out the complete list of piece sizes in an optimal decomposition of a rod of length  $n$ :

PRINT-CUT-ROD-SOLUTION( $p, n$ )

```

1  ( $r, s$ ) = EXTENDED-BOTTOM-UP-CUT-ROD( $p, n$ )
2  while  $n > 0$ 
3      print  $s[n]$ 
4       $n = n - s[n]$ 
```

In our rod-cutting example, the call EXTENDED-BOTTOM-UP-CUT-ROD( $p, 10$ ) would return the following arrays:

$i$	0	1	2	3	4	5	6	7	8	9	10
$r[i]$	0	1	5	8	10	13	17	18	22	25	30
$s[i]$		1	2	3	2	2	6	1	2	3	10

A call to PRINT-CUT-ROD-SOLUTION( $p, 10$ ) would print just 10, but a call with  $n = 7$  would print the cuts 1 and 6, corresponding to the first optimal decomposition for  $r_7$  given earlier.

---

## 15.2 Matrix-chain multiplication

Our next example of dynamic programming is an algorithm that solves the problem of matrix-chain multiplication. We are given a sequence (chain)  $\langle A_1, A_2, \dots, A_n \rangle$  of  $n$  matrices to be multiplied, and we wish to compute the product

$$A_1 A_2 \cdots A_n . \tag{15.5}$$

We can evaluate the expression (15.5) using the standard algorithm for multiplying pairs of matrices as a subroutine once we have parenthesized it to resolve all ambiguities in how the matrices are multiplied together. Matrix multiplication is associative, and so all parenthesizations yield the same product. A product of matrices is **fully parenthesized** if it is either a single matrix or the product of two fully parenthesized matrix products, surrounded by parentheses. For example, if the chain of matrices is  $\langle A_1, A_2, A_3, A_4 \rangle$ , then we can fully parenthesize the product  $A_1 A_2 A_3 A_4$  in five distinct ways:

$(A_1(A_2(A_3A_4)))$  ,  
 $(A_1((A_2A_3)A_4))$  ,  
 $((A_1A_2)(A_3A_4))$  ,  
 $((A_1(A_2A_3))A_4)$  ,  
 $((A_1A_2)A_3)A_4$  .

How we parenthesize a chain of matrices can have a dramatic impact on the cost of evaluating the product. Consider first the cost of multiplying two matrices. The standard algorithm is given by the following pseudocode, which generalizes the SQUARE-MATRIX-MULTIPLY procedure from Section 4.2. The attributes *rows* and *columns* are the numbers of rows and columns in a matrix.

```

MATRIX-MULTIPLY(A, B)
1  if A.columns ≠ B.rows
2      error “incompatible dimensions”
3  else let C be a new A.rows × B.columns matrix
4      for i = 1 to A.rows
5          for j = 1 to B.columns
6              cij = 0
7              for k = 1 to A.columns
8                  cij = cij + aik · bkj
9  return C
  
```

We can multiply two matrices *A* and *B* only if they are **compatible**: the number of columns of *A* must equal the number of rows of *B*. If *A* is a  $p \times q$  matrix and *B* is a  $q \times r$  matrix, the resulting matrix *C* is a  $p \times r$  matrix. The time to compute *C* is dominated by the number of scalar multiplications in line 8, which is  $pqr$ . In what follows, we shall express costs in terms of the number of scalar multiplications.

To illustrate the different costs incurred by different parenthesizations of a matrix product, consider the problem of a chain  $\langle A_1, A_2, A_3 \rangle$  of three matrices. Suppose that the dimensions of the matrices are  $10 \times 100$ ,  $100 \times 5$ , and  $5 \times 50$ , respectively. If we multiply according to the parenthesization  $((A_1A_2)A_3)$ , we perform  $10 \cdot 100 \cdot 5 = 5000$  scalar multiplications to compute the  $10 \times 5$  matrix product  $A_1A_2$ , plus another  $10 \cdot 5 \cdot 50 = 2500$  scalar multiplications to multiply this matrix by  $A_3$ , for a total of 7500 scalar multiplications. If instead we multiply according to the parenthesization  $(A_1(A_2A_3))$ , we perform  $100 \cdot 5 \cdot 50 = 25,000$  scalar multiplications to compute the  $100 \times 50$  matrix product  $A_2A_3$ , plus another  $10 \cdot 100 \cdot 50 = 50,000$  scalar multiplications to multiply  $A_1$  by this matrix, for a total of 75,000 scalar multiplications. Thus, computing the product according to the first parenthesization is 10 times faster.

We state the **matrix-chain multiplication problem** as follows: given a chain  $\langle A_1, A_2, \dots, A_n \rangle$  of  $n$  matrices, where for  $i = 1, 2, \dots, n$ , matrix  $A_i$  has dimension

$p_{i-1} \times p_i$ , fully parenthesize the product  $A_1 A_2 \cdots A_n$  in a way that minimizes the number of scalar multiplications.

Note that in the matrix-chain multiplication problem, we are not actually multiplying matrices. Our goal is only to determine an order for multiplying matrices that has the lowest cost. Typically, the time invested in determining this optimal order is more than paid for by the time saved later on when actually performing the matrix multiplications (such as performing only 7500 scalar multiplications instead of 75,000).

### Counting the number of parenthesizations

Before solving the matrix-chain multiplication problem by dynamic programming, let us convince ourselves that exhaustively checking all possible parenthesizations does not yield an efficient algorithm. Denote the number of alternative parenthesizations of a sequence of  $n$  matrices by  $P(n)$ . When  $n = 1$ , we have just one matrix and therefore only one way to fully parenthesize the matrix product. When  $n \geq 2$ , a fully parenthesized matrix product is the product of two fully parenthesized matrix subproducts, and the split between the two subproducts may occur between the  $k$ th and  $(k + 1)$ st matrices for any  $k = 1, 2, \dots, n - 1$ . Thus, we obtain the recurrence

$$P(n) = \begin{cases} 1 & \text{if } n = 1, \\ \sum_{k=1}^{n-1} P(k)P(n-k) & \text{if } n \geq 2. \end{cases} \quad (15.6)$$

Problem 12-4 asked you to show that the solution to a similar recurrence is the sequence of **Catalan numbers**, which grows as  $\Omega(4^n/n^{3/2})$ . A simpler exercise (see Exercise 15.2-3) is to show that the solution to the recurrence (15.6) is  $\Omega(2^n)$ . The number of solutions is thus exponential in  $n$ , and the brute-force method of exhaustive search makes for a poor strategy when determining how to optimally parenthesize a matrix chain.

### Applying dynamic programming

We shall use the dynamic-programming method to determine how to optimally parenthesize a matrix chain. In so doing, we shall follow the four-step sequence that we stated at the beginning of this chapter:

1. Characterize the structure of an optimal solution.
2. Recursively define the value of an optimal solution.
3. Compute the value of an optimal solution.

#### 4. Construct an optimal solution from computed information.

We shall go through these steps in order, demonstrating clearly how we apply each step to the problem.

### Step 1: The structure of an optimal parenthesization

For our first step in the dynamic-programming paradigm, we find the optimal substructure and then use it to construct an optimal solution to the problem from optimal solutions to subproblems. In the matrix-chain multiplication problem, we can perform this step as follows. For convenience, let us adopt the notation  $A_{i..j}$ , where  $i \leq j$ , for the matrix that results from evaluating the product  $A_i A_{i+1} \cdots A_j$ . Observe that if the problem is nontrivial, i.e.,  $i < j$ , then to parenthesize the product  $A_i A_{i+1} \cdots A_j$ , we must split the product between  $A_k$  and  $A_{k+1}$  for some integer  $k$  in the range  $i \leq k < j$ . That is, for some value of  $k$ , we first compute the matrices  $A_{i..k}$  and  $A_{k+1..j}$  and then multiply them together to produce the final product  $A_{i..j}$ . The cost of parenthesizing this way is the cost of computing the matrix  $A_{i..k}$ , plus the cost of computing  $A_{k+1..j}$ , plus the cost of multiplying them together.

The optimal substructure of this problem is as follows. Suppose that to optimally parenthesize  $A_i A_{i+1} \cdots A_j$ , we split the product between  $A_k$  and  $A_{k+1}$ . Then the way we parenthesize the “prefix” subchain  $A_i A_{i+1} \cdots A_k$  within this optimal parenthesization of  $A_i A_{i+1} \cdots A_j$  must be an optimal parenthesization of  $A_i A_{i+1} \cdots A_k$ . Why? If there were a less costly way to parenthesize  $A_i A_{i+1} \cdots A_k$ , then we could substitute that parenthesization in the optimal parenthesization of  $A_i A_{i+1} \cdots A_j$  to produce another way to parenthesize  $A_i A_{i+1} \cdots A_j$  whose cost was lower than the optimum: a contradiction. A similar observation holds for how we parenthesize the subchain  $A_{k+1} A_{k+2} \cdots A_j$  in the optimal parenthesization of  $A_i A_{i+1} \cdots A_j$ : it must be an optimal parenthesization of  $A_{k+1} A_{k+2} \cdots A_j$ .

Now we use our optimal substructure to show that we can construct an optimal solution to the problem from optimal solutions to subproblems. We have seen that any solution to a nontrivial instance of the matrix-chain multiplication problem requires us to split the product, and that any optimal solution contains within it optimal solutions to subproblem instances. Thus, we can build an optimal solution to an instance of the matrix-chain multiplication problem by splitting the problem into two subproblems (optimally parenthesizing  $A_i A_{i+1} \cdots A_k$  and  $A_{k+1} A_{k+2} \cdots A_j$ ), finding optimal solutions to subproblem instances, and then combining these optimal subproblem solutions. We must ensure that when we search for the correct place to split the product, we have considered all possible places, so that we are sure of having examined the optimal one.

### Step 2: A recursive solution

Next, we define the cost of an optimal solution recursively in terms of the optimal solutions to subproblems. For the matrix-chain multiplication problem, we pick as our subproblems the problems of determining the minimum cost of parenthesizing  $A_i A_{i+1} \cdots A_j$  for  $1 \leq i \leq j \leq n$ . Let  $m[i, j]$  be the minimum number of scalar multiplications needed to compute the matrix  $A_{i..j}$ ; for the full problem, the lowest-cost way to compute  $A_{1..n}$  would thus be  $m[1, n]$ .

We can define  $m[i, j]$  recursively as follows. If  $i = j$ , the problem is trivial; the chain consists of just one matrix  $A_{i..i} = A_i$ , so that no scalar multiplications are necessary to compute the product. Thus,  $m[i, i] = 0$  for  $i = 1, 2, \dots, n$ . To compute  $m[i, j]$  when  $i < j$ , we take advantage of the structure of an optimal solution from step 1. Let us assume that to optimally parenthesize, we split the product  $A_i A_{i+1} \cdots A_j$  between  $A_k$  and  $A_{k+1}$ , where  $i \leq k < j$ . Then,  $m[i, j]$  equals the minimum cost for computing the subproducts  $A_{i..k}$  and  $A_{k+1..j}$ , plus the cost of multiplying these two matrices together. Recalling that each matrix  $A_i$  is  $p_{i-1} \times p_i$ , we see that computing the matrix product  $A_{i..k} A_{k+1..j}$  takes  $p_{i-1} p_k p_j$  scalar multiplications. Thus, we obtain

$$m[i, j] = m[i, k] + m[k + 1, j] + p_{i-1} p_k p_j .$$

This recursive equation assumes that we know the value of  $k$ , which we do not. There are only  $j - i$  possible values for  $k$ , however, namely  $k = i, i + 1, \dots, j - 1$ . Since the optimal parenthesization must use one of these values for  $k$ , we need only check them all to find the best. Thus, our recursive definition for the minimum cost of parenthesizing the product  $A_i A_{i+1} \cdots A_j$  becomes

$$m[i, j] = \begin{cases} 0 & \text{if } i = j , \\ \min_{i \leq k < j} \{m[i, k] + m[k + 1, j] + p_{i-1} p_k p_j\} & \text{if } i < j . \end{cases} \quad (15.7)$$

The  $m[i, j]$  values give the costs of optimal solutions to subproblems, but they do not provide all the information we need to construct an optimal solution. To help us do so, we define  $s[i, j]$  to be a value of  $k$  at which we split the product  $A_i A_{i+1} \cdots A_j$  in an optimal parenthesization. That is,  $s[i, j]$  equals a value  $k$  such that  $m[i, j] = m[i, k] + m[k + 1, j] + p_{i-1} p_k p_j$ .

### Step 3: Computing the optimal costs

At this point, we could easily write a recursive algorithm based on recurrence (15.7) to compute the minimum cost  $m[1, n]$  for multiplying  $A_1 A_2 \cdots A_n$ . As we saw for the rod-cutting problem, and as we shall see in Section 15.3, this recursive algorithm takes exponential time, which is no better than the brute-force method of checking each way of parenthesizing the product.

Observe that we have relatively few distinct subproblems: one subproblem for each choice of  $i$  and  $j$  satisfying  $1 \leq i \leq j \leq n$ , or  $\binom{n}{2} + n = \Theta(n^2)$  in all. A recursive algorithm may encounter each subproblem many times in different branches of its recursion tree. This property of overlapping subproblems is the second hallmark of when dynamic programming applies (the first hallmark being optimal substructure).

Instead of computing the solution to recurrence (15.7) recursively, we compute the optimal cost by using a tabular, bottom-up approach. (We present the corresponding top-down approach using memoization in Section 15.3.)

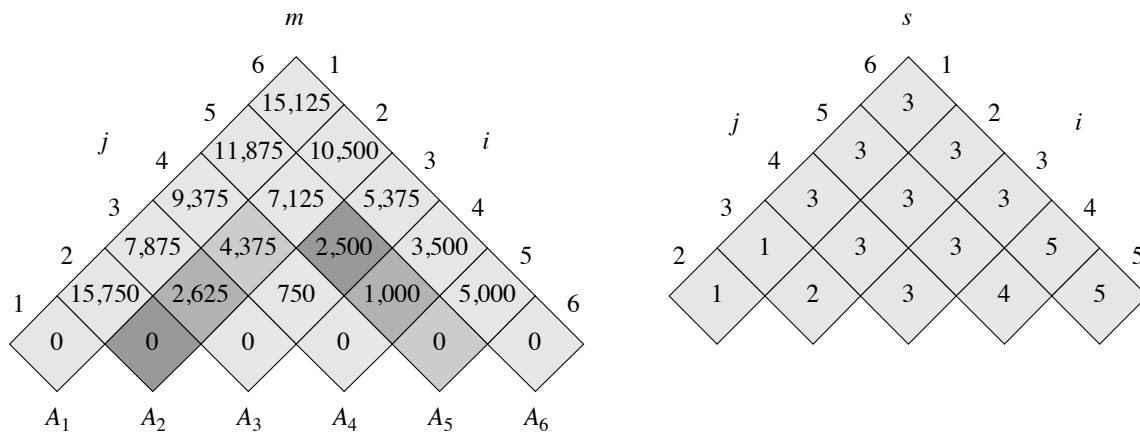
We shall implement the tabular, bottom-up method in the procedure MATRIX-CHAIN-ORDER, which appears below. This procedure assumes that matrix  $A_i$  has dimensions  $p_{i-1} \times p_i$  for  $i = 1, 2, \dots, n$ . Its input is a sequence  $p = \langle p_0, p_1, \dots, p_n \rangle$ , where  $p.length = n + 1$ . The procedure uses an auxiliary table  $m[1..n, 1..n]$  for storing the  $m[i, j]$  costs and another auxiliary table  $s[1..n-1, 2..n]$  that records which index of  $k$  achieved the optimal cost in computing  $m[i, j]$ . We shall use the table  $s$  to construct an optimal solution.

In order to implement the bottom-up approach, we must determine which entries of the table we refer to when computing  $m[i, j]$ . Equation (15.7) shows that the cost  $m[i, j]$  of computing a matrix-chain product of  $j - i + 1$  matrices depends only on the costs of computing matrix-chain products of fewer than  $j - i + 1$  matrices. That is, for  $k = i, i + 1, \dots, j - 1$ , the matrix  $A_{i..k}$  is a product of  $k - i + 1 < j - i + 1$  matrices and the matrix  $A_{k+1..j}$  is a product of  $j - k < j - i + 1$  matrices. Thus, the algorithm should fill in the table  $m$  in a manner that corresponds to solving the parenthesization problem on matrix chains of increasing length. For the subproblem of optimally parenthesizing the chain  $A_i A_{i+1} \dots A_j$ , we consider the subproblem size to be the length  $j - i + 1$  of the chain.

MATRIX-CHAIN-ORDER( $p$ )

```

1   $n = p.length - 1$ 
2  let  $m[1..n, 1..n]$  and  $s[1..n-1, 2..n]$  be new tables
3  for  $i = 1$  to  $n$ 
4       $m[i, i] = 0$ 
5  for  $l = 2$  to  $n$            //  $l$  is the chain length
6      for  $i = 1$  to  $n - l + 1$ 
7           $j = i + l - 1$ 
8           $m[i, j] = \infty$ 
9          for  $k = i$  to  $j - 1$ 
10              $q = m[i, k] + m[k + 1, j] + p_{i-1} p_k p_j$ 
11             if  $q < m[i, j]$ 
12                  $m[i, j] = q$ 
13                  $s[i, j] = k$ 
14  return  $m$  and  $s$ 
```



**Figure 15.5** The  $m$  and  $s$  tables computed by MATRIX-CHAIN-ORDER for  $n = 6$  and the following matrix dimensions:

matrix	$A_1$	$A_2$	$A_3$	$A_4$	$A_5$	$A_6$
dimension	$30 \times 35$	$35 \times 15$	$15 \times 5$	$5 \times 10$	$10 \times 20$	$20 \times 25$

The tables are rotated so that the main diagonal runs horizontally. The  $m$  table uses only the main diagonal and upper triangle, and the  $s$  table uses only the upper triangle. The minimum number of scalar multiplications to multiply the 6 matrices is  $m[1, 6] = 15,125$ . Of the darker entries, the pairs that have the same shading are taken together in line 10 when computing

$$\begin{aligned}
 m[2, 5] &= \min \begin{cases} m[2, 2] + m[3, 5] + p_1 p_2 p_5 = 0 + 2500 + 35 \cdot 15 \cdot 20 = 13,000, \\ m[2, 3] + m[4, 5] + p_1 p_3 p_5 = 2625 + 1000 + 35 \cdot 5 \cdot 20 = 7125, \\ m[2, 4] + m[5, 5] + p_1 p_4 p_5 = 4375 + 0 + 35 \cdot 10 \cdot 20 = 11,375 \end{cases} \\
 &= 7125.
 \end{aligned}$$

The algorithm first computes  $m[i, i] = 0$  for  $i = 1, 2, \dots, n$  (the minimum costs for chains of length 1) in lines 3–4. It then uses recurrence (15.7) to compute  $m[i, i + 1]$  for  $i = 1, 2, \dots, n - 1$  (the minimum costs for chains of length  $l = 2$ ) during the first execution of the **for** loop in lines 5–13. The second time through the loop, it computes  $m[i, i + 2]$  for  $i = 1, 2, \dots, n - 2$  (the minimum costs for chains of length  $l = 3$ ), and so forth. At each step, the  $m[i, j]$  cost computed in lines 10–13 depends only on table entries  $m[i, k]$  and  $m[k + 1, j]$  already computed.

Figure 15.5 illustrates this procedure on a chain of  $n = 6$  matrices. Since we have defined  $m[i, j]$  only for  $i \leq j$ , only the portion of the table  $m$  on or above the main diagonal is used. The figure shows the table rotated to make the main diagonal run horizontally. The matrix chain is listed along the bottom. Using this layout, we can find the minimum cost  $m[i, j]$  for multiplying a subchain  $A_i A_{i+1} \cdots A_j$  of matrices at the intersection of lines running northeast from  $A_i$  and northwest



from  $A_j$ . Each horizontal row in the table contains the entries for matrix chains of the same length. MATRIX-CHAIN-ORDER computes the rows from bottom to top and from left to right within each row. It computes each entry  $m[i, j]$  using the products  $p_{i-1}p_kp_j$  for  $k = i, i + 1, \dots, j - 1$  and all entries southwest and southeast from  $m[i, j]$ .

A simple inspection of the nested loop structure of MATRIX-CHAIN-ORDER yields a running time of  $O(n^3)$  for the algorithm. The loops are nested three deep, and each loop index ( $l, i$ , and  $k$ ) takes on at most  $n - 1$  values. Exercise 15.2-5 asks you to show that the running time of this algorithm is in fact also  $\Omega(n^3)$ . The algorithm requires  $\Theta(n^2)$  space to store the  $m$  and  $s$  tables. Thus, MATRIX-CHAIN-ORDER is much more efficient than the exponential-time method of enumerating all possible parenthesizations and checking each one.

#### Step 4: Constructing an optimal solution

Although MATRIX-CHAIN-ORDER determines the optimal number of scalar multiplications needed to compute a matrix-chain product, it does not directly show how to multiply the matrices. The table  $s[1..n - 1, 2..n]$  gives us the information we need to do so. Each entry  $s[i, j]$  records a value of  $k$  such that an optimal parenthesization of  $A_i A_{i+1} \dots A_j$  splits the product between  $A_k$  and  $A_{k+1}$ . Thus, we know that the final matrix multiplication in computing  $A_{1..n}$  optimally is  $A_{1..s[1,n]}A_{s[1,n]+1..n}$ . We can determine the earlier matrix multiplications recursively, since  $s[1, s[1, n]]$  determines the last matrix multiplication when computing  $A_{1..s[1,n]}$  and  $s[s[1, n] + 1, n]$  determines the last matrix multiplication when computing  $A_{s[1,n]+1..n}$ . The following recursive procedure prints an optimal parenthesization of  $\langle A_i, A_{i+1}, \dots, A_j \rangle$ , given the  $s$  table computed by MATRIX-CHAIN-ORDER and the indices  $i$  and  $j$ . The initial call PRINT-OPTIMAL-PARENS( $s, 1, n$ ) prints an optimal parenthesization of  $\langle A_1, A_2, \dots, A_n \rangle$ .

```

PRINT-OPTIMAL-PARENS( $s, i, j$ )
1  if  $i == j$ 
2      print " $A$ " $i$ 
3  else print "("
4      PRINT-OPTIMAL-PARENS( $s, i, s[i, j]$ )
5      PRINT-OPTIMAL-PARENS( $s, s[i, j] + 1, j$ )
6      print ")"

```

In the example of Figure 15.5, the call PRINT-OPTIMAL-PARENS( $s, 1, 6$ ) prints the parenthesization  $((A_1(A_2A_3))((A_4A_5)A_6))$ .

---

### 15.3 Elements of dynamic programming

Although we have just worked through two examples of the dynamic-programming method, you might still be wondering just when the method applies. From an engineering perspective, when should we look for a dynamic-programming solution to a problem? In this section, we examine the two key ingredients that an opti-

mization problem must have in order for dynamic programming to apply: optimal substructure and overlapping subproblems. We also revisit and discuss more fully how memoization might help us take advantage of the overlapping-subproblems property in a top-down recursive approach.

### Optimal substructure

The first step in solving an optimization problem by dynamic programming is to characterize the structure of an optimal solution. Recall that a problem exhibits **optimal substructure** if an optimal solution to the problem contains within it optimal solutions to subproblems. Whenever a problem exhibits optimal substructure, we have a good clue that dynamic programming might apply. (As Chapter 16 discusses, it also might mean that a greedy strategy applies, however.) In dynamic programming, we build an optimal solution to the problem from optimal solutions to subproblems. Consequently, we must take care to ensure that the range of subproblems we consider includes those used in an optimal solution.

We discovered optimal substructure in both of the problems we have examined in this chapter so far. In Section 15.1, we observed that the optimal way of cutting up a rod of length  $n$  (if we make any cuts at all) involves optimally cutting up the two pieces resulting from the first cut. In Section 15.2, we observed that an optimal parenthesization of  $A_i A_{i+1} \cdots A_j$  that splits the product between  $A_k$  and  $A_{k+1}$  contains within it optimal solutions to the problems of parenthesizing  $A_i A_{i+1} \cdots A_k$  and  $A_{k+1} A_{k+2} \cdots A_j$ .

You will find yourself following a common pattern in discovering optimal substructure:

1. You show that a solution to the problem consists of making a choice, such as choosing an initial cut in a rod or choosing an index at which to split the matrix chain. Making this choice leaves one or more subproblems to be solved.
2. You suppose that for a given problem, you are given the choice that leads to an optimal solution. You do not concern yourself yet with how to determine this choice. You just assume that it has been given to you.
3. Given this choice, you determine which subproblems ensue and how to best characterize the resulting space of subproblems.
4. You show that the solutions to the subproblems used within an optimal solution to the problem must themselves be optimal by using a “cut-and-paste” technique. You do so by supposing that each of the subproblem solutions is not optimal and then deriving a contradiction. In particular, by “cutting out” the nonoptimal solution to each subproblem and “pasting in” the optimal one, you show that you can get a better solution to the original problem, thus contradicting your supposition that you already had an optimal solution. If an optimal

solution gives rise to more than one subproblem, they are typically so similar that you can modify the cut-and-paste argument for one to apply to the others with little effort.

To characterize the space of subproblems, a good rule of thumb says to try to keep the space as simple as possible and then expand it as necessary. For example, the space of subproblems that we considered for the rod-cutting problem contained the problems of optimally cutting up a rod of length  $i$  for each size  $i$ . This subproblem space worked well, and we had no need to try a more general space of subproblems.

Conversely, suppose that we had tried to constrain our subproblem space for matrix-chain multiplication to matrix products of the form  $A_1 A_2 \cdots A_j$ . As before, an optimal parenthesization must split this product between  $A_k$  and  $A_{k+1}$  for some  $1 \leq k < j$ . Unless we could guarantee that  $k$  always equals  $j - 1$ , we would find that we had subproblems of the form  $A_1 A_2 \cdots A_k$  and  $A_{k+1} A_{k+2} \cdots A_j$ , and that the latter subproblem is not of the form  $A_1 A_2 \cdots A_j$ . For this problem, we needed to allow our subproblems to vary at “both ends,” that is, to allow both  $i$  and  $j$  to vary in the subproblem  $A_i A_{i+1} \cdots A_j$ .

Optimal substructure varies across problem domains in two ways:

1. how many subproblems an optimal solution to the original problem uses, and
2. how many choices we have in determining which subproblem(s) to use in an optimal solution.

In the rod-cutting problem, an optimal solution for cutting up a rod of size  $n$  uses just one subproblem (of size  $n - i$ ), but we must consider  $n$  choices for  $i$  in order to determine which one yields an optimal solution. Matrix-chain multiplication for the subchain  $A_i A_{i+1} \cdots A_j$  serves as an example with two subproblems and  $j - i$  choices. For a given matrix  $A_k$  at which we split the product, we have two subproblems—parenthesizing  $A_i A_{i+1} \cdots A_k$  and parenthesizing  $A_{k+1} A_{k+2} \cdots A_j$ —and we must solve *both* of them optimally. Once we determine the optimal solutions to subproblems, we choose from among  $j - i$  candidates for the index  $k$ .

Informally, the running time of a dynamic-programming algorithm depends on the product of two factors: the number of subproblems overall and how many choices we look at for each subproblem. In rod cutting, we had  $\Theta(n)$  subproblems overall, and at most  $n$  choices to examine for each, yielding an  $O(n^2)$  running time. Matrix-chain multiplication had  $\Theta(n^2)$  subproblems overall, and in each we had at most  $n - 1$  choices, giving an  $O(n^3)$  running time (actually, a  $\Theta(n^3)$  running time, by Exercise 15.2-5).

Usually, the subproblem graph gives an alternative way to perform the same analysis. Each vertex corresponds to a subproblem, and the choices for a sub-

problem are the edges incident from that subproblem. Recall that in rod cutting, the subproblem graph had  $n$  vertices and at most  $n$  edges per vertex, yielding an  $O(n^2)$  running time. For matrix-chain multiplication, if we were to draw the subproblem graph, it would have  $\Theta(n^2)$  vertices and each vertex would have degree at most  $n - 1$ , giving a total of  $O(n^3)$  vertices and edges.

Dynamic programming often uses optimal substructure in a bottom-up fashion. That is, we first find optimal solutions to subproblems and, having solved the subproblems, we find an optimal solution to the problem. Finding an optimal solution to the problem entails making a choice among subproblems as to which we will use in solving the problem. The cost of the problem solution is usually the subproblem costs plus a cost that is directly attributable to the choice itself. In rod cutting, for example, first we solved the subproblems of determining optimal ways to cut up rods of length  $i$  for  $i = 0, 1, \dots, n - 1$ , and then we determined which such subproblem yielded an optimal solution for a rod of length  $n$ , using equation (15.2). The cost attributable to the choice itself is the term  $p_i$  in equation (15.2). In matrix-chain multiplication, we determined optimal parenthesizations of subchains of  $A_i A_{i+1} \cdots A_j$ , and then we chose the matrix  $A_k$  at which to split the product. The cost attributable to the choice itself is the term  $p_{i-1} p_k p_j$ .

In Chapter 16, we shall examine “greedy algorithms,” which have many similarities to dynamic programming. In particular, problems to which greedy algorithms apply have optimal substructure. One major difference between greedy algorithms and dynamic programming is that instead of first finding optimal solutions to subproblems and then making an informed choice, greedy algorithms first make a “greedy” choice—the choice that looks best at the time—and then solve a resulting subproblem, without bothering to solve all possible related smaller subproblems. Surprisingly, in some cases this strategy works!

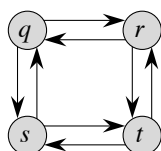
### Subtleties

You should be careful not to assume that optimal substructure applies when it does not. Consider the following two problems in which we are given a directed graph  $G = (V, E)$  and vertices  $u, v \in V$ .

**Unweighted shortest path:**<sup>3</sup> Find a path from  $u$  to  $v$  consisting of the fewest edges. Such a path must be simple, since removing a cycle from a path produces a path with fewer edges.

---

<sup>3</sup>We use the term “unweighted” to distinguish this problem from that of finding shortest paths with weighted edges, which we shall see in Chapters 24 and 25. We can use the breadth-first search technique of Chapter 22 to solve the unweighted problem.



**Figure 15.6** A directed graph showing that the problem of finding a longest simple path in an unweighted directed graph does not have optimal substructure. The path  $q \rightarrow r \rightarrow t$  is a longest simple path from  $q$  to  $t$ , but the subpath  $q \rightarrow r$  is not a longest simple path from  $q$  to  $r$ , nor is the subpath  $r \rightarrow t$  a longest simple path from  $r$  to  $t$ .

**Unweighted longest simple path:** Find a simple path from  $u$  to  $v$  consisting of the most edges. We need to include the requirement of simplicity because otherwise we can traverse a cycle as many times as we like to create paths with an arbitrarily large number of edges.

The unweighted shortest-path problem exhibits optimal substructure, as follows. Suppose that  $u \neq v$ , so that the problem is nontrivial. Then, any path  $p$  from  $u$  to  $v$  must contain an intermediate vertex, say  $w$ . (Note that  $w$  may be  $u$  or  $v$ .) Thus, we can decompose the path  $u \xrightarrow{p} v$  into subpaths  $u \xrightarrow{p_1} w \xrightarrow{p_2} v$ . Clearly, the number of edges in  $p$  equals the number of edges in  $p_1$  plus the number of edges in  $p_2$ . We claim that if  $p$  is an optimal (i.e., shortest) path from  $u$  to  $v$ , then  $p_1$  must be a shortest path from  $u$  to  $w$ . Why? We use a “cut-and-paste” argument: if there were another path, say  $p'_1$ , from  $u$  to  $w$  with fewer edges than  $p_1$ , then we could cut out  $p_1$  and paste in  $p'_1$  to produce a path  $u \xrightarrow{p'_1} w \xrightarrow{p_2} v$  with fewer edges than  $p$ , thus contradicting  $p$ ’s optimality. Symmetrically,  $p_2$  must be a shortest path from  $w$  to  $v$ . Thus, we can find a shortest path from  $u$  to  $v$  by considering all intermediate vertices  $w$ , finding a shortest path from  $u$  to  $w$  and a shortest path from  $w$  to  $v$ , and choosing an intermediate vertex  $w$  that yields the overall shortest path. In Section 25.2, we use a variant of this observation of optimal substructure to find a shortest path between every pair of vertices on a weighted, directed graph.

You might be tempted to assume that the problem of finding an unweighted longest simple path exhibits optimal substructure as well. After all, if we decompose a longest simple path  $u \xrightarrow{p} v$  into subpaths  $u \xrightarrow{p_1} w \xrightarrow{p_2} v$ , then mustn’t  $p_1$  be a longest simple path from  $u$  to  $w$ , and mustn’t  $p_2$  be a longest simple path from  $w$  to  $v$ ? The answer is no! Figure 15.6 supplies an example. Consider the path  $q \rightarrow r \rightarrow t$ , which is a longest simple path from  $q$  to  $t$ . Is  $q \rightarrow r$  a longest simple path from  $q$  to  $r$ ? No, for the path  $q \rightarrow s \rightarrow t \rightarrow r$  is a simple path that is longer. Is  $r \rightarrow t$  a longest simple path from  $r$  to  $t$ ? No again, for the path  $r \rightarrow q \rightarrow s \rightarrow t$  is a simple path that is longer.

This example shows that for longest simple paths, not only does the problem lack optimal substructure, but we cannot necessarily assemble a “legal” solution to the problem from solutions to subproblems. If we combine the longest simple paths  $q \rightarrow s \rightarrow t \rightarrow r$  and  $r \rightarrow q \rightarrow s \rightarrow t$ , we get the path  $q \rightarrow s \rightarrow t \rightarrow r \rightarrow q \rightarrow s \rightarrow t$ , which is not simple. Indeed, the problem of finding an unweighted longest simple path does not appear to have any sort of optimal substructure. No efficient dynamic-programming algorithm for this problem has ever been found. In fact, this problem is NP-complete, which—as we shall see in Chapter 34—means that we are unlikely to find a way to solve it in polynomial time.

Why is the substructure of a longest simple path so different from that of a shortest path? Although a solution to a problem for both longest and shortest paths uses two subproblems, the subproblems in finding the longest simple path are not *independent*, whereas for shortest paths they are. What do we mean by subproblems being independent? We mean that the solution to one subproblem does not affect the solution to another subproblem of the same problem. For the example of Figure 15.6, we have the problem of finding a longest simple path from  $q$  to  $t$  with two subproblems: finding longest simple paths from  $q$  to  $r$  and from  $r$  to  $t$ . For the first of these subproblems, we choose the path  $q \rightarrow s \rightarrow t \rightarrow r$ , and so we have also used the vertices  $s$  and  $t$ . We can no longer use these vertices in the second subproblem, since the combination of the two solutions to subproblems would yield a path that is not simple. If we cannot use vertex  $t$  in the second problem, then we cannot solve it at all, since  $t$  is required to be on the path that we find, and it is not the vertex at which we are “splicing” together the subproblem solutions (that vertex being  $r$ ). Because we use vertices  $s$  and  $t$  in one subproblem solution, we cannot use them in the other subproblem solution. We must use at least one of them to solve the other subproblem, however, and we must use both of them to solve it optimally. Thus, we say that these subproblems are not independent. Looked at another way, using resources in solving one subproblem (those resources being vertices) renders them unavailable for the other subproblem.

Why, then, are the subproblems independent for finding a shortest path? The answer is that by nature, the subproblems do not share resources. We claim that if a vertex  $w$  is on a shortest path  $p$  from  $u$  to  $v$ , then we can splice together *any* shortest path  $u \xrightarrow{p_1} w$  and *any* shortest path  $w \xrightarrow{p_2} v$  to produce a shortest path from  $u$  to  $v$ . We are assured that, other than  $w$ , no vertex can appear in both paths  $p_1$  and  $p_2$ . Why? Suppose that some vertex  $x \neq w$  appears in both  $p_1$  and  $p_2$ , so that we can decompose  $p_1$  as  $u \xrightarrow{p_{ux}} x \rightsquigarrow w$  and  $p_2$  as  $w \rightsquigarrow x \xrightarrow{p_{xv}} v$ . By the optimal substructure of this problem, path  $p$  has as many edges as  $p_1$  and  $p_2$  together; let’s say that  $p$  has  $e$  edges. Now let us construct a path  $p' = u \xrightarrow{p_{ux}} x \xrightarrow{p_{xv}} v$  from  $u$  to  $v$ . Because we have excised the paths from  $x$  to  $w$  and from  $w$  to  $x$ , each of which contains at least one edge, path  $p'$  contains at most  $e - 2$  edges, which contradicts

the assumption that  $p$  is a shortest path. Thus, we are assured that the subproblems for the shortest-path problem are independent.

Both problems examined in Sections 15.1 and 15.2 have independent subproblems. In matrix-chain multiplication, the subproblems are multiplying subchains  $A_i A_{i+1} \cdots A_k$  and  $A_{k+1} A_{k+2} \cdots A_j$ . These subchains are disjoint, so that no matrix could possibly be included in both of them. In rod cutting, to determine the best way to cut up a rod of length  $n$ , we look at the best ways of cutting up rods of length  $i$  for  $i = 0, 1, \dots, n-1$ . Because an optimal solution to the length- $n$  problem includes just one of these subproblem solutions (after we have cut off the first piece), independence of subproblems is not an issue.

### Overlapping subproblems

The second ingredient that an optimization problem must have for dynamic programming to apply is that the space of subproblems must be “small” in the sense that a recursive algorithm for the problem solves the same subproblems over and over, rather than always generating new subproblems. Typically, the total number of distinct subproblems is a polynomial in the input size. When a recursive algorithm revisits the same problem repeatedly, we say that the optimization problem has *overlapping subproblems*.<sup>4</sup> In contrast, a problem for which a divide-and-conquer approach is suitable usually generates brand-new problems at each step of the recursion. Dynamic-programming algorithms typically take advantage of overlapping subproblems by solving each subproblem once and then storing the solution in a table where it can be looked up when needed, using constant time per lookup.

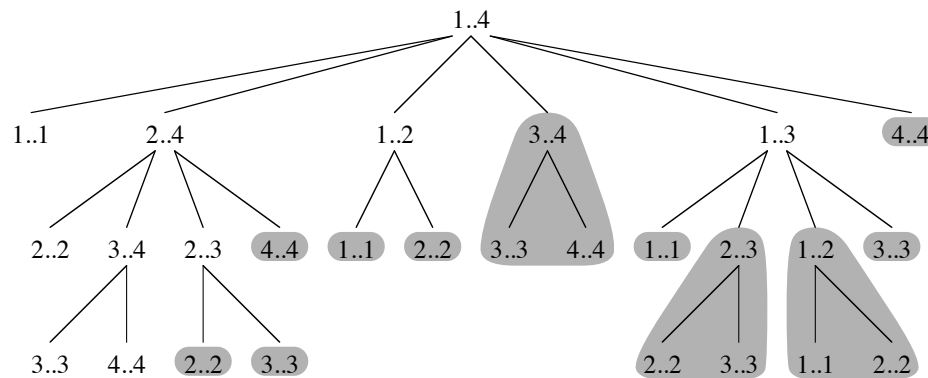
In Section 15.1, we briefly examined how a recursive solution to rod cutting makes exponentially many calls to find solutions of smaller subproblems. Our dynamic-programming solution takes an exponential-time recursive algorithm down to quadratic time.

To illustrate the overlapping-subproblems property in greater detail, let us re-examine the matrix-chain multiplication problem. Referring back to Figure 15.5, observe that MATRIX-CHAIN-ORDER repeatedly looks up the solution to subproblems in lower rows when solving subproblems in higher rows. For example, it references entry  $m[3, 4]$  four times: during the computations of  $m[2, 4]$ ,  $m[1, 4]$ ,

---

<sup>4</sup>It may seem strange that dynamic programming relies on subproblems being both independent and overlapping. Although these requirements may sound contradictory, they describe two different notions, rather than two points on the same axis. Two subproblems of the same problem are independent if they do not share resources. Two subproblems are overlapping if they are really the same subproblem that occurs as a subproblem of different problems.





$m[3, 5]$ , and  $m[3, 6]$ . If we were to recompute  $m[3, 4]$  each time, rather than just looking it up, the running time would increase dramatically. To see how, consider the following (inefficient) recursive procedure that determines  $m[i, j]$ , the minimum number of scalar multiplications needed to compute the matrix-chain product  $A_{i..j} = A_i A_{i+1} \cdots A_j$ . The procedure is based directly on the recurrence (15.7).

```

1  if  $i == j$ 
2      return 0
3   $m[i, j] = \infty$ 
4  for  $k = i$  to  $j - 1$ 
5       $q = \text{RECURSIVE-MATRIX-CHAIN}(p, i, k)$ 
         $+ \text{RECURSIVE-MATRIX-CHAIN}(p, k + 1, j)$ 
         $+ p_{i-1} p_k p_j$ 
6      if  $q < m[i, j]$ 
7           $m[i, j] = q$ 
8  return  $m[i, j]$ 

```

In fact, we can show that the time to compute  $m[1, n]$  by this recursive procedure is at least exponential in  $n$ . Let  $T(n)$  denote the time taken by RECURSIVE-MATRIX-CHAIN to compute an optimal parenthesization of a chain of  $n$  matrices. Because the execution of lines 1–2 and of lines 6–7 each take at least unit time, as

does the multiplication in line 5, inspection of the procedure yields the recurrence

$$\begin{aligned} T(1) &\geq 1, \\ T(n) &\geq 1 + \sum_{k=1}^{n-1} (T(k) + T(n-k) + 1) \quad \text{for } n > 1. \end{aligned}$$

Noting that for  $i = 1, 2, \dots, n-1$ , each term  $T(i)$  appears once as  $T(k)$  and once as  $T(n-k)$ , and collecting the  $n-1$  1s in the summation together with the 1 out front, we can rewrite the recurrence as

$$T(n) \geq 2 \sum_{i=1}^{n-1} T(i) + n. \quad (15.8)$$

We shall prove that  $T(n) = \Omega(2^n)$  using the substitution method. Specifically, we shall show that  $T(n) \geq 2^{n-1}$  for all  $n \geq 1$ . The basis is easy, since  $T(1) \geq 1 = 2^0$ . Inductively, for  $n \geq 2$  we have

$$\begin{aligned} T(n) &\geq 2 \sum_{i=1}^{n-1} 2^{i-1} + n \\ &= 2 \sum_{i=0}^{n-2} 2^i + n \\ &= 2(2^{n-1} - 1) + n \quad (\text{by equation (A.5)}) \\ &= 2^n - 2 + n \\ &\geq 2^{n-1}, \end{aligned}$$

which completes the proof. Thus, the total amount of work performed by the call `RECURSIVE-MATRIX-CHAIN( $p, 1, n$ )` is at least exponential in  $n$ .

Compare this top-down, recursive algorithm (without memoization) with the bottom-up dynamic-programming algorithm. The latter is more efficient because it takes advantage of the overlapping-subproblems property. Matrix-chain multiplication has only  $\Theta(n^2)$  distinct subproblems, and the dynamic-programming algorithm solves each exactly once. The recursive algorithm, on the other hand, must again solve each subproblem every time it reappears in the recursion tree. Whenever a recursion tree for the natural recursive solution to a problem contains the same subproblem repeatedly, and the total number of distinct subproblems is small, dynamic programming can improve efficiency, sometimes dramatically.

### Reconstructing an optimal solution

As a practical matter, we often store which choice we made in each subproblem in a table so that we do not have to reconstruct this information from the costs that we stored.

For matrix-chain multiplication, the table  $s[i, j]$  saves us a significant amount of work when reconstructing an optimal solution. Suppose that we did not maintain the  $s[i, j]$  table, having filled in only the table  $m[i, j]$  containing optimal subproblem costs. We choose from among  $j - i$  possibilities when we determine which subproblems to use in an optimal solution to parenthesizing  $A_i A_{i+1} \cdots A_j$ , and  $j - i$  is not a constant. Therefore, it would take  $\Theta(j - i) = \omega(1)$  time to reconstruct which subproblems we chose for a solution to a given problem. By storing in  $s[i, j]$  the index of the matrix at which we split the product  $A_i A_{i+1} \cdots A_j$ , we can reconstruct each choice in  $O(1)$  time.

### Memoization

As we saw for the rod-cutting problem, there is an alternative approach to dynamic programming that often offers the efficiency of the bottom-up dynamic-programming approach while maintaining a top-down strategy. The idea is to **memoize** the natural, but inefficient, recursive algorithm. As in the bottom-up approach, we maintain a table with subproblem solutions, but the control structure for filling in the table is more like the recursive algorithm.

A memoized recursive algorithm maintains an entry in a table for the solution to each subproblem. Each table entry initially contains a special value to indicate that the entry has yet to be filled in. When the subproblem is first encountered as the recursive algorithm unfolds, its solution is computed and then stored in the table. Each subsequent time that we encounter this subproblem, we simply look up the value stored in the table and return it.<sup>5</sup>

Here is a memoized version of RECURSIVE-MATRIX-CHAIN. Note where it resembles the memoized top-down method for the rod-cutting problem.

---

<sup>5</sup>This approach presupposes that we know the set of all possible subproblem parameters and that we have established the relationship between table positions and subproblems. Another, more general, approach is to memoize by using hashing with the subproblem parameters as keys.

MEMOIZED-MATRIX-CHAIN( $p$ )

```

1   $n = p.length - 1$ 
2  let  $m[1..n, 1..n]$  be a new table
3  for  $i = 1$  to  $n$ 
4      for  $j = i$  to  $n$ 
5           $m[i, j] = \infty$ 
6  return LOOKUP-CHAIN( $m, p, 1, n$ )

```

LOOKUP-CHAIN( $m, p, i, j$ )

```

1  if  $m[i, j] < \infty$ 
2      return  $m[i, j]$ 
3  if  $i == j$ 
4       $m[i, j] = 0$ 
5  else for  $k = i$  to  $j - 1$ 
6       $q = \text{LOOKUP-CHAIN}(m, p, i, k)$ 
           $+ \text{LOOKUP-CHAIN}(m, p, k + 1, j) + p_{i-1}p_kp_j$ 
7      if  $q < m[i, j]$ 
8           $m[i, j] = q$ 
9  return  $m[i, j]$ 

```

The MEMOIZED-MATRIX-CHAIN procedure, like MATRIX-CHAIN-ORDER, maintains a table  $m[1..n, 1..n]$  of computed values of  $m[i, j]$ , the minimum number of scalar multiplications needed to compute the matrix  $A_{i..j}$ . Each table entry initially contains the value  $\infty$  to indicate that the entry has yet to be filled in. Upon calling LOOKUP-CHAIN( $m, p, i, j$ ), if line 1 finds that  $m[i, j] < \infty$ , then the procedure simply returns the previously computed cost  $m[i, j]$  in line 2. Otherwise, the cost is computed as in RECURSIVE-MATRIX-CHAIN, stored in  $m[i, j]$ , and returned. Thus, LOOKUP-CHAIN( $m, p, i, j$ ) always returns the value of  $m[i, j]$ , but it computes it only upon the first call of LOOKUP-CHAIN with these specific values of  $i$  and  $j$ .

Figure 15.7 illustrates how MEMOIZED-MATRIX-CHAIN saves time compared with RECURSIVE-MATRIX-CHAIN. Shaded subtrees represent values that it looks up rather than recomputes.

Like the bottom-up dynamic-programming algorithm MATRIX-CHAIN-ORDER, the procedure MEMOIZED-MATRIX-CHAIN runs in  $O(n^3)$  time. Line 5 of MEMOIZED-MATRIX-CHAIN executes  $\Theta(n^2)$  times. We can categorize the calls of LOOKUP-CHAIN into two types:

1. calls in which  $m[i, j] = \infty$ , so that lines 3–9 execute, and
2. calls in which  $m[i, j] < \infty$ , so that LOOKUP-CHAIN simply returns in line 2.

There are  $\Theta(n^2)$  calls of the first type, one per table entry. All calls of the second type are made as recursive calls by calls of the first type. Whenever a given call of LOOKUP-CHAIN makes recursive calls, it makes  $O(n)$  of them. Therefore, there are  $O(n^3)$  calls of the second type in all. Each call of the second type takes  $O(1)$  time, and each call of the first type takes  $O(n)$  time plus the time spent in its recursive calls. The total time, therefore, is  $O(n^3)$ . Memoization thus turns an  $\Omega(2^n)$ -time algorithm into an  $O(n^3)$ -time algorithm.

In summary, we can solve the matrix-chain multiplication problem by either a top-down, memoized dynamic-programming algorithm or a bottom-up dynamic-programming algorithm in  $O(n^3)$  time. Both methods take advantage of the overlapping-subproblems property. There are only  $\Theta(n^2)$  distinct subproblems in total, and either of these methods computes the solution to each subproblem only once. Without memoization, the natural recursive algorithm runs in exponential time, since solved subproblems are repeatedly solved.

In general practice, if all subproblems must be solved at least once, a bottom-up dynamic-programming algorithm usually outperforms the corresponding top-down memoized algorithm by a constant factor, because the bottom-up algorithm has no overhead for recursion and less overhead for maintaining the table. Moreover, for some problems we can exploit the regular pattern of table accesses in the dynamic-programming algorithm to reduce time or space requirements even further. Alternatively, if some subproblems in the subproblem space need not be solved at all, the memoized solution has the advantage of solving only those subproblems that are definitely required.

In this chapter, we consider the problem of finding shortest paths between all pairs of vertices in a graph. This problem might arise in making a table of distances between all pairs of cities for a road atlas. As in Chapter 24, we are given a weighted, directed graph  $G = (V, E)$  with a weight function  $w : E \rightarrow \mathbf{R}$  that maps edges to real-valued weights. We wish to find, for every pair of vertices  $u, v \in V$ , a shortest (least-weight) path from  $u$  to  $v$ , where the weight of a path is the sum of the weights of its constituent edges. We typically want the output in tabular form: the entry in  $u$ 's row and  $v$ 's column should be the weight of a shortest path from  $u$  to  $v$ .

We can solve an all-pairs shortest-paths problem by running a single-source shortest-paths algorithm  $|V|$  times, once for each vertex as the source. If all edge weights are nonnegative, we can use Dijkstra's algorithm. If we use the linear-array implementation of the min-priority queue, the running time is  $O(V^3 + VE) = O(V^3)$ . The binary min-heap implementation of the min-priority queue yields a running time of  $O(VE \lg V)$ , which is an improvement if the graph is sparse. Alternatively, we can implement the min-priority queue with a Fibonacci heap, yielding a running time of  $O(V^2 \lg V + VE)$ .

If negative-weight edges are allowed, Dijkstra's algorithm can no longer be used. Instead, we must run the slower Bellman-Ford algorithm once from each vertex. The resulting running time is  $O(V^2E)$ , which on a dense graph is  $O(V^4)$ . In this chapter we shall see how to do better. We shall also investigate the relation of the all-pairs shortest-paths problem to matrix multiplication and study its algebraic structure.

Unlike the single-source algorithms, which assume an adjacency-list representation of the graph, most of the algorithms in this chapter use an adjacency-matrix representation. (Johnson's algorithm for sparse graphs uses adjacency lists.) For convenience, we assume that the vertices are numbered  $1, 2, \dots, |V|$ , so that the input is an  $n \times n$  matrix  $W$  representing the edge weights of an  $n$ -vertex directed graph  $G = (V, E)$ . That is,  $W = (w_{ij})$ , where

$$w_{ij} = \begin{cases} 0 & \text{if } i = j, \\ \text{the weight of directed edge } (i, j) & \text{if } i \neq j \text{ and } (i, j) \in E, \\ \infty & \text{if } i \neq j \text{ and } (i, j) \notin E. \end{cases} \quad (25.1)$$

Negative-weight edges are allowed, but we assume for the time being that the input graph contains no negative-weight cycles.

The tabular output of the all-pairs shortest-paths algorithms presented in this chapter is an  $n \times n$  matrix  $D = (d_{ij})$ , where entry  $d_{ij}$  contains the weight of a shortest path from vertex  $i$  to vertex  $j$ . That is, if we let  $\delta(i, j)$  denote the shortest-path weight from vertex  $i$  to vertex  $j$  (as in Chapter 24), then  $d_{ij} = \delta(i, j)$  at termination.

To solve the all-pairs shortest-paths problem on an input adjacency matrix, we need to compute not only the shortest-path weights but also a **predecessor matrix**  $\Pi = (\pi_{ij})$ , where  $\pi_{ij}$  is NIL if either  $i = j$  or there is no path from  $i$  to  $j$ , and otherwise  $\pi_{ij}$  is the predecessor of  $j$  on some shortest path from  $i$ . Just as the predecessor subgraph  $G_\pi$  from Chapter 24 is a shortest-paths tree for a given source vertex, the subgraph induced by the  $i$ th row of the  $\Pi$  matrix should be a shortest-paths tree with root  $i$ . For each vertex  $i \in V$ , we define the **predecessor subgraph** of  $G$  for  $i$  as  $G_{\pi,i} = (V_{\pi,i}, E_{\pi,i})$ , where

$$V_{\pi,i} = \{j \in V : \pi_{ij} \neq \text{NIL}\} \cup \{i\}$$

and

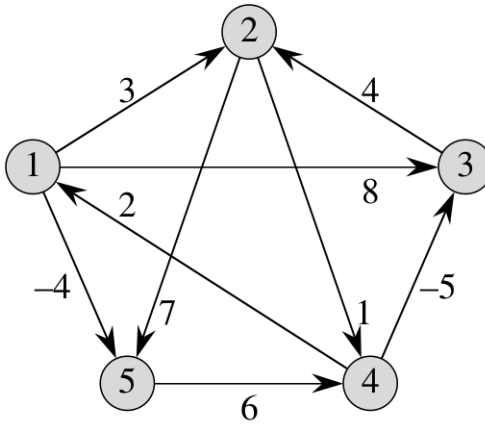
$$E_{\pi,i} = \{(\pi_{ij}, j) : j \in V_{\pi,i} - \{i\}\}.$$

If  $G_{\pi,i}$  is a shortest-paths tree, then the following procedure, which is a modified version of the PRINT-PATH procedure from Chapter 22, prints a shortest path from vertex  $i$  to vertex  $j$ .

PRINT-ALL-PAIRS-SHORTEST-PATH( $\Pi, i, j$ )

```

1  if  $i = j$ 
2    then print  $i$ 
3  else if  $\pi_{ij} = \text{NIL}$ 
4    then print “no path from”  $i$  “to”  $j$  “exists”
5    else PRINT-ALL-PAIRS-SHORTEST-PATH( $\Pi, i, \pi_{ij}$ )
6    print  $j$ 
```



An example graph that is referred to as Figure 25.1 in the following content on the Floyd-Warshall algorithm.



## 25.2 The Floyd-Warshall algorithm

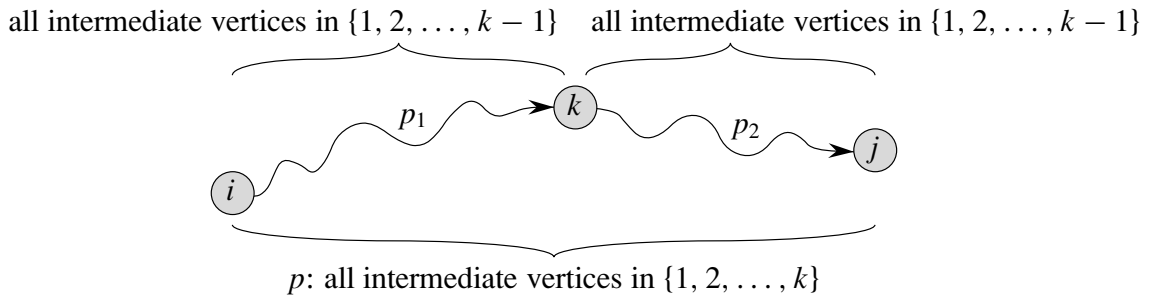
In this section, we shall use a different dynamic-programming formulation to solve the all-pairs shortest-paths problem on a directed graph  $G = (V, E)$ . The resulting algorithm, known as the **Floyd-Warshall algorithm**, runs in  $\Theta(V^3)$  time. As before, negative-weight edges may be present, but we assume that there are no negative-weight cycles. As in Section 25.1, we shall follow the dynamic-programming process to develop the algorithm. After studying the resulting algorithm, we shall present a similar method for finding the transitive closure of a directed graph.

### The structure of a shortest path

In the Floyd-Warshall algorithm, we use a different characterization of the structure of a shortest path than we used in the matrix-multiplication-based all-pairs algorithms. The algorithm considers the “intermediate” vertices of a shortest path, where an *intermediate* vertex of a simple path  $p = \langle v_1, v_2, \dots, v_l \rangle$  is any vertex of  $p$  other than  $v_1$  or  $v_l$ , that is, any vertex in the set  $\{v_2, v_3, \dots, v_{l-1}\}$ .

The Floyd-Warshall algorithm is based on the following observation. Under our assumption that the vertices of  $G$  are  $V = \{1, 2, \dots, n\}$ , let us consider a subset  $\{1, 2, \dots, k\}$  of vertices for some  $k$ . For any pair of vertices  $i, j \in V$ , consider all paths from  $i$  to  $j$  whose intermediate vertices are all drawn from  $\{1, 2, \dots, k\}$ , and let  $p$  be a minimum-weight path from among them. (Path  $p$  is simple.) The Floyd-Warshall algorithm exploits a relationship between path  $p$  and shortest paths from  $i$  to  $j$  with all intermediate vertices in the set  $\{1, 2, \dots, k-1\}$ . The relationship depends on whether or not  $k$  is an intermediate vertex of path  $p$ .

- If  $k$  is not an intermediate vertex of path  $p$ , then all intermediate vertices of path  $p$  are in the set  $\{1, 2, \dots, k-1\}$ . Thus, a shortest path from vertex  $i$  to vertex  $j$  with all intermediate vertices in the set  $\{1, 2, \dots, k-1\}$  is also a shortest path from  $i$  to  $j$  with all intermediate vertices in the set  $\{1, 2, \dots, k\}$ .
- If  $k$  is an intermediate vertex of path  $p$ , then we break  $p$  down into  $i \xrightarrow{p_1} k \xrightarrow{p_2} j$  as shown in Figure 25.3. By Lemma 24.1,  $p_1$  is a shortest path from  $i$  to  $k$  with all intermediate vertices in the set  $\{1, 2, \dots, k\}$ . Because vertex  $k$  is not an intermediate vertex of path  $p_1$ , we see that  $p_1$  is a shortest path from  $i$  to  $k$  with all intermediate vertices in the set  $\{1, 2, \dots, k-1\}$ . Similarly,  $p_2$  is a shortest path from vertex  $k$  to vertex  $j$  with all intermediate vertices in the set  $\{1, 2, \dots, k-1\}$ .



**Figure 25.3** Path  $p$  is a shortest path from vertex  $i$  to vertex  $j$ , and  $k$  is the highest-numbered intermediate vertex of  $p$ . Path  $p_1$ , the portion of path  $p$  from vertex  $i$  to vertex  $k$ , has all intermediate vertices in the set  $\{1, 2, \dots, k-1\}$ . The same holds for path  $p_2$  from vertex  $k$  to vertex  $j$ .

### A recursive solution to the all-pairs shortest-paths problem

Based on the above observations, we define a recursive formulation of shortest-path estimates that is different from the one in Section 25.1. Let  $d_{ij}^{(k)}$  be the weight of a shortest path from vertex  $i$  to vertex  $j$  for which all intermediate vertices are in the set  $\{1, 2, \dots, k\}$ . When  $k = 0$ , a path from vertex  $i$  to vertex  $j$  with no intermediate vertex numbered higher than 0 has no intermediate vertices at all. Such a path has at most one edge, and hence  $d_{ij}^{(0)} = w_{ij}$ . A recursive definition following the above discussion is given by

$$d_{ij}^{(k)} = \begin{cases} w_{ij} & \text{if } k = 0, \\ \min(d_{ij}^{(k-1)}, d_{ik}^{(k-1)} + d_{kj}^{(k-1)}) & \text{if } k \geq 1. \end{cases} \quad (25.5)$$

Because for any path, all intermediate vertices are in the set  $\{1, 2, \dots, n\}$ , the matrix  $D^{(n)} = (d_{ij}^{(n)})$  gives the final answer:  $d_{ij}^{(n)} = \delta(i, j)$  for all  $i, j \in V$ .

### Computing the shortest-path weights bottom up

Based on recurrence (25.5), the following bottom-up procedure can be used to compute the values  $d_{ij}^{(k)}$  in order of increasing values of  $k$ . Its input is an  $n \times n$  matrix  $W$  defined as in equation (25.1). The procedure returns the matrix  $D^{(n)}$  of shortest-path weights.

FLOYD-WARSHALL( $W$ )

```

1   $n \leftarrow \text{rows}[W]$ 
2   $D^{(0)} \leftarrow W$ 
3  for  $k \leftarrow 1$  to  $n$ 
4      do for  $i \leftarrow 1$  to  $n$ 
5          do for  $j \leftarrow 1$  to  $n$ 
6              do  $d_{ij}^{(k)} \leftarrow \min(d_{ij}^{(k-1)}, d_{ik}^{(k-1)} + d_{kj}^{(k-1)})$ 
7  return  $D^{(n)}$ 
```

$$\begin{aligned}
D^{(0)} &= \begin{pmatrix} 0 & 3 & 8 & \infty & -4 \\ \infty & 0 & \infty & 1 & 7 \\ \infty & 4 & 0 & \infty & \infty \\ 2 & \infty & -5 & 0 & \infty \\ \infty & \infty & \infty & 6 & 0 \end{pmatrix} & \Pi^{(0)} &= \begin{pmatrix} \text{NIL} & 1 & 1 & \text{NIL} & 1 \\ \text{NIL} & \text{NIL} & \text{NIL} & 2 & 2 \\ \text{NIL} & 3 & \text{NIL} & \text{NIL} & \text{NIL} \\ 4 & \text{NIL} & 4 & \text{NIL} & \text{NIL} \\ \text{NIL} & \text{NIL} & \text{NIL} & 5 & \text{NIL} \end{pmatrix} \\
D^{(1)} &= \begin{pmatrix} 0 & 3 & 8 & \infty & -4 \\ \infty & 0 & \infty & 1 & 7 \\ \infty & 4 & 0 & \infty & \infty \\ 2 & 5 & -5 & 0 & -2 \\ \infty & \infty & \infty & 6 & 0 \end{pmatrix} & \Pi^{(1)} &= \begin{pmatrix} \text{NIL} & 1 & 1 & \text{NIL} & 1 \\ \text{NIL} & \text{NIL} & \text{NIL} & 2 & 2 \\ \text{NIL} & 3 & \text{NIL} & \text{NIL} & \text{NIL} \\ 4 & 1 & 4 & \text{NIL} & 1 \\ \text{NIL} & \text{NIL} & \text{NIL} & 5 & \text{NIL} \end{pmatrix} \\
D^{(2)} &= \begin{pmatrix} 0 & 3 & 8 & 4 & -4 \\ \infty & 0 & \infty & 1 & 7 \\ \infty & 4 & 0 & 5 & 11 \\ 2 & 5 & -5 & 0 & -2 \\ \infty & \infty & \infty & 6 & 0 \end{pmatrix} & \Pi^{(2)} &= \begin{pmatrix} \text{NIL} & 1 & 1 & 2 & 1 \\ \text{NIL} & \text{NIL} & \text{NIL} & 2 & 2 \\ \text{NIL} & 3 & \text{NIL} & 2 & 2 \\ 4 & 1 & 4 & \text{NIL} & 1 \\ \text{NIL} & \text{NIL} & \text{NIL} & 5 & \text{NIL} \end{pmatrix} \\
D^{(3)} &= \begin{pmatrix} 0 & 3 & 8 & 4 & -4 \\ \infty & 0 & \infty & 1 & 7 \\ \infty & 4 & 0 & 5 & 11 \\ 2 & -1 & -5 & 0 & -2 \\ \infty & \infty & \infty & 6 & 0 \end{pmatrix} & \Pi^{(3)} &= \begin{pmatrix} \text{NIL} & 1 & 1 & 2 & 1 \\ \text{NIL} & \text{NIL} & \text{NIL} & 2 & 2 \\ \text{NIL} & 3 & \text{NIL} & 2 & 2 \\ 4 & 3 & 4 & \text{NIL} & 1 \\ \text{NIL} & \text{NIL} & \text{NIL} & 5 & \text{NIL} \end{pmatrix} \\
D^{(4)} &= \begin{pmatrix} 0 & 3 & -1 & 4 & -4 \\ 3 & 0 & -4 & 1 & -1 \\ 7 & 4 & 0 & 5 & 3 \\ 2 & -1 & -5 & 0 & -2 \\ 8 & 5 & 1 & 6 & 0 \end{pmatrix} & \Pi^{(4)} &= \begin{pmatrix} \text{NIL} & 1 & 4 & 2 & 1 \\ 4 & \text{NIL} & 4 & 2 & 1 \\ 4 & 3 & \text{NIL} & 2 & 1 \\ 4 & 3 & 4 & \text{NIL} & 1 \\ 4 & 3 & 4 & 5 & \text{NIL} \end{pmatrix} \\
D^{(5)} &= \begin{pmatrix} 0 & 1 & -3 & 2 & -4 \\ 3 & 0 & -4 & 1 & -1 \\ 7 & 4 & 0 & 5 & 3 \\ 2 & -1 & -5 & 0 & -2 \\ 8 & 5 & 1 & 6 & 0 \end{pmatrix} & \Pi^{(5)} &= \begin{pmatrix} \text{NIL} & 3 & 4 & 5 & 1 \\ 4 & \text{NIL} & 4 & 2 & 1 \\ 4 & 3 & \text{NIL} & 2 & 1 \\ 4 & 3 & 4 & \text{NIL} & 1 \\ 4 & 3 & 4 & 5 & \text{NIL} \end{pmatrix}
\end{aligned}$$

**Figure 25.4** The sequence of matrices  $D^{(k)}$  and  $\Pi^{(k)}$  computed by the Floyd-Warshall algorithm for the graph in Figure 25.1.

Figure 25.4 shows the matrices  $D^{(k)}$  computed by the Floyd-Warshall algorithm for the graph in Figure 25.1.

The running time of the Floyd-Warshall algorithm is determined by the triply nested **for** loops of lines 3–6. Because each execution of line 6 takes  $O(1)$  time, the algorithm runs in time  $\Theta(n^3)$ . As in the final algorithm in Section 25.1, the

code is tight, with no elaborate data structures, and so the constant hidden in the  $\Theta$ -notation is small. Thus, the Floyd-Warshall algorithm is quite practical for even moderate-sized input graphs.

### Constructing a shortest path

There are a variety of different methods for constructing shortest paths in the Floyd-Warshall algorithm. One way is to compute the matrix  $D$  of shortest-path weights and then construct the predecessor matrix  $\Pi$  from the  $D$  matrix. This method can be implemented to run in  $O(n^3)$  time (Exercise 25.1-6). Given the predecessor matrix  $\Pi$ , the PRINT-ALL-PAIRS-SHORTEST-PATH procedure can be used to print the vertices on a given shortest path.

We can compute the predecessor matrix  $\Pi$  “on-line” just as the Floyd-Warshall algorithm computes the matrices  $D^{(k)}$ . Specifically, we compute a sequence of matrices  $\Pi^{(0)}, \Pi^{(1)}, \dots, \Pi^{(n)}$ , where  $\Pi = \Pi^{(n)}$  and  $\pi_{ij}^{(k)}$  is defined to be the predecessor of vertex  $j$  on a shortest path from vertex  $i$  with all intermediate vertices in the set  $\{1, 2, \dots, k\}$ .

We can give a recursive formulation of  $\pi_{ij}^{(k)}$ . When  $k = 0$ , a shortest path from  $i$  to  $j$  has no intermediate vertices at all. Thus,

$$\pi_{ij}^{(0)} = \begin{cases} \text{NIL} & \text{if } i = j \text{ or } w_{ij} = \infty, \\ i & \text{if } i \neq j \text{ and } w_{ij} < \infty. \end{cases} \quad (25.6)$$

For  $k \geq 1$ , if we take the path  $i \rightsquigarrow k \rightsquigarrow j$ , where  $k \neq j$ , then the predecessor of  $j$  we choose is the same as the predecessor of  $j$  we chose on a shortest path from  $k$  with all intermediate vertices in the set  $\{1, 2, \dots, k-1\}$ . Otherwise, we choose the same predecessor of  $j$  that we chose on a shortest path from  $i$  with all intermediate vertices in the set  $\{1, 2, \dots, k-1\}$ . Formally, for  $k \geq 1$ ,

$$\pi_{ij}^{(k)} = \begin{cases} \pi_{ij}^{(k-1)} & \text{if } d_{ij}^{(k-1)} \leq d_{ik}^{(k-1)} + d_{kj}^{(k-1)}, \\ \pi_{kj}^{(k-1)} & \text{if } d_{ij}^{(k-1)} > d_{ik}^{(k-1)} + d_{kj}^{(k-1)}. \end{cases} \quad (25.7)$$

We leave the incorporation of the  $\Pi^{(k)}$  matrix computations into the FLOYD-WARSHALL procedure as Exercise 25.2-3. Figure 25.4 shows the sequence of  $\Pi^{(k)}$  matrices that the resulting algorithm computes for the graph of Figure 25.1. The exercise also asks for the more difficult task of proving that the predecessor subgraph  $G_{\pi,i}$  is a shortest-paths tree with root  $i$ . Yet another way to reconstruct shortest paths is given as Exercise 25.2-7.