

8

Dynamic Programming

An idea, like a ghost . . . must be spoken to a little before it will explain itself.

—Charles Dickens (1812–1870)

Dynamic programming is an algorithm design technique with a rather interesting history. It was invented by a prominent U.S. mathematician, Richard Bellman, in the 1950s as a general method for optimizing multistage decision processes. Thus, the word “programming” in the name of this technique stands for “planning” and does not refer to computer programming. After proving its worth as an important tool of applied mathematics, dynamic programming has eventually come to be considered, at least in computer science circles, as a general algorithm design technique that does not have to be limited to special types of optimization problems. It is from this point of view that we will consider this technique here.

Dynamic programming is a technique for solving problems with overlapping subproblems. Typically, these subproblems arise from a recurrence relating a given problem’s solution to solutions of its smaller subproblems. Rather than solving overlapping subproblems again and again, dynamic programming suggests solving each of the smaller subproblems only once and recording the results in a table from which a solution to the original problem can then be obtained.

This technique can be illustrated by revisiting the Fibonacci numbers discussed in Section 2.5. (If you have not read that section, you will be able to follow the discussion anyway. But it is a beautiful topic, so if you feel a temptation to read it, do succumb to it.) The Fibonacci numbers are the elements of the sequence

0, 1, 1, 2, 3, 5, 8, 13, 21, 34, . . . ,

which can be defined by the simple recurrence

$$F(n) = F(n - 1) + F(n - 2) \quad \text{for } n > 1 \quad (8.1)$$

and two initial conditions

$$F(0) = 0, \quad F(1) = 1. \quad (8.2)$$

If we try to use recurrence (8.1) directly to compute the n th Fibonacci number $F(n)$, we would have to recompute the same values of this function many times (see Figure 2.6 for an example). Note that the problem of computing $F(n)$ is expressed in terms of its smaller and overlapping subproblems of computing $F(n-1)$ and $F(n-2)$. So we can simply fill elements of a one-dimensional array with the $n+1$ consecutive values of $F(n)$ by starting, in view of initial conditions (8.2), with 0 and 1 and using equation (8.1) as the rule for producing all the other elements. Obviously, the last element of this array will contain $F(n)$. Single-loop pseudocode of this very simple algorithm can be found in Section 2.5.

Note that we can, in fact, avoid using an extra array to accomplish this task by recording the values of just the last two elements of the Fibonacci sequence (Problem 8 in Exercises 2.5). This phenomenon is not unusual, and we shall encounter it in a few more examples in this chapter. Thus, although a straightforward application of dynamic programming can be interpreted as a special variety of space-for-time trade-off, a dynamic programming algorithm can sometimes be refined to avoid using extra space.

Certain algorithms compute the n th Fibonacci number without computing all the preceding elements of this sequence (see Section 2.5). It is typical of an algorithm based on the classic bottom-up dynamic programming approach, however, to solve *all* smaller subproblems of a given problem. One variation of the dynamic programming approach seeks to avoid solving unnecessary subproblems. This technique, illustrated in Section 8.2, exploits so-called memory functions and can be considered a top-down variation of dynamic programming.

Whether one uses the classical bottom-up version of dynamic programming or its top-down variation, the crucial step in designing such an algorithm remains the same: deriving a recurrence relating a solution to the problem to solutions to its smaller subproblems. The immediate availability of equation (8.1) for computing the n th Fibonacci number is one of the few exceptions to this rule.

Since a majority of dynamic programming applications deal with optimization problems, we also need to mention a general principle that underlines such applications. Richard Bellman called it the ***principle of optimality***. In terms somewhat different from its original formulation, it says that an optimal solution to any instance of an optimization problem is composed of optimal solutions to its subinstances. The principle of optimality holds much more often than not. (To give a rather rare example, it fails for finding the longest simple path in a graph.) Although its applicability to a particular problem needs to be checked, of course, such a check is usually not a principal difficulty in developing a dynamic programming algorithm.

In the sections and exercises of this chapter are a few standard examples of dynamic programming algorithms. (The algorithms in Section 8.4 were, in fact,

invented independently of the discovery of dynamic programming and only later came to be viewed as examples of this technique's applications.) Numerous other applications range from the optimal way of breaking text into lines (e.g., [Baa00]) to image resizing [Avi07] to a variety of applications to sophisticated engineering problems (e.g., [Ber01]).

8.1 Three Basic Examples

The goal of this section is to introduce dynamic programming via three typical examples.

EXAMPLE 1 *Coin-row problem* There is a row of n coins whose values are some positive integers c_1, c_2, \dots, c_n , not necessarily distinct. The goal is to pick up the maximum amount of money subject to the constraint that no two coins adjacent in the initial row can be picked up.

Let $F(n)$ be the maximum amount that can be picked up from the row of n coins. To derive a recurrence for $F(n)$, we partition all the allowed coin selections into two groups: those that include the last coin and those without it. The largest amount we can get from the first group is equal to $c_n + F(n - 2)$ —the value of the n th coin plus the maximum amount we can pick up from the first $n - 2$ coins. The maximum amount we can get from the second group is equal to $F(n - 1)$ by the definition of $F(n)$. Thus, we have the following recurrence subject to the obvious initial conditions:

$$\begin{aligned} F(n) &= \max\{c_n + F(n - 2), F(n - 1)\} \quad \text{for } n > 1, \\ F(0) &= 0, \quad F(1) = c_1. \end{aligned} \tag{8.3}$$

We can compute $F(n)$ by filling the one-row table left to right in the manner similar to the way it was done for the n th Fibonacci number by Algorithm *Fib*(n) in Section 2.5.

ALGORITHM *CoinRow*($C[1..n]$)

```
//Applies formula (8.3) bottom up to find the maximum amount of money
//that can be picked up from a coin row without picking two adjacent coins
//Input: Array  $C[1..n]$  of positive integers indicating the coin values
//Output: The maximum amount of money that can be picked up
 $F[0] \leftarrow 0; \quad F[1] \leftarrow C[1]$ 
for  $i \leftarrow 2$  to  $n$  do
     $F[i] \leftarrow \max(C[i] + F[i - 2], F[i - 1])$ 
return  $F[n]$ 
```

The application of the algorithm to the coin row of denominations 5, 1, 2, 10, 6, 2 is shown in Figure 8.1. It yields the maximum amount of 17. It is worth pointing

		index	0	1	2	3	4	5	6
		C		5	1	2	10	6	2
$F[0] = 0, F[1] = c_1 = 5$		F	0	5					

		index	0	1	2	3	4	5	6
		C		5	1	2	10	6	2
$F[2] = \max\{1 + 0, 5\} = 5$		F	0	5	5				

		index	0	1	2	3	4	5	6
		C		5	1	2	10	6	2
$F[3] = \max\{2 + 5, 5\} = 7$		F	0	5	5	7			

		index	0	1	2	3	4	5	6
		C		5	1	2	10	6	2
$F[4] = \max\{10 + 5, 7\} = 15$		F	0	5	5	7	15		

		index	0	1	2	3	4	5	6
		C		5	1	2	10	6	2
$F[5] = \max\{6 + 7, 15\} = 15$		F	0	5	5	7	15	15	

		index	0	1	2	3	4	5	6
		C		5	1	2	10	6	2
$F[6] = \max\{2 + 15, 15\} = 17$		F	0	5	5	7	15	15	17

FIGURE 8.1 Solving the coin-row problem by dynamic programming for the coin row 5, 1, 2, 10, 6, 2.

out that, in fact, we also solved the problem for the first i coins in the row given for every $1 \leq i \leq 6$. For example, for $i = 3$, the maximum amount is $F(3) = 7$.

To find the coins with the maximum total value found, we need to backtrace the computations to see which of the two possibilities— $c_n + F(n - 2)$ or $F(n - 1)$ —produced the maxima in formula (8.3). In the last application of the formula, it was the sum $c_6 + F(4)$, which means that the coin $c_6 = 2$ is a part of an optimal solution. Moving to computing $F(4)$, the maximum was produced by the sum $c_4 + F(2)$, which means that the coin $c_4 = 10$ is a part of an optimal solution as well. Finally, the maximum in computing $F(2)$ was produced by $F(1)$, implying that the coin c_2 is not the part of an optimal solution and the coin $c_1 = 5$ is. Thus, the optimal solution is $\{c_1, c_4, c_6\}$. To avoid repeating the same computations during the backtracing, the information about which of the two terms in (8.3) was larger can be recorded in an extra array when the values of F are computed.

Using the *CoinRow* to find $F(n)$, the largest amount of money that can be picked up, as well as the coins composing an optimal set, clearly takes $\Theta(n)$ time and $\Theta(n)$ space. This is by far superior to the alternatives: the straightforward top-

down application of recurrence (8.3) and solving the problem by exhaustive search (Problem 3 in this section's exercises). ■

EXAMPLE 2 *Change-making problem* Consider the general instance of the following well-known problem. Give change for amount n using the minimum number of coins of denominations $d_1 < d_2 < \dots < d_m$. For the coin denominations used in the United States, as for those used in most if not all other countries, there is a very simple and efficient algorithm discussed in the next chapter. Here, we consider a dynamic programming algorithm for the general case, assuming availability of unlimited quantities of coins for each of the m denominations $d_1 < d_2 < \dots < d_m$ where $d_1 = 1$.

Let $F(n)$ be the minimum number of coins whose values add up to n ; it is convenient to define $F(0) = 0$. The amount n can only be obtained by adding one coin of denomination d_j to the amount $n - d_j$ for $j = 1, 2, \dots, m$ such that $n \geq d_j$. Therefore, we can consider all such denominations and select the one minimizing $F(n - d_j) + 1$. Since 1 is a constant, we can, of course, find the smallest $F(n - d_j)$ first and then add 1 to it. Hence, we have the following recurrence for $F(n)$:

$$\begin{aligned} F(n) &= \min_{j: n \geq d_j} \{F(n - d_j)\} + 1 \quad \text{for } n > 0, \\ F(0) &= 0. \end{aligned} \tag{8.4}$$

We can compute $F(n)$ by filling a one-row table left to right in the manner similar to the way it was done above for the coin-row problem, but computing a table entry here requires finding the minimum of up to m numbers.

ALGORITHM *ChangeMaking*($D[1..m], n$)

```
//Applies dynamic programming to find the minimum number of coins
//of denominations  $d_1 < d_2 < \dots < d_m$  where  $d_1 = 1$  that add up to a
//given amount  $n$ 
//Input: Positive integer  $n$  and array  $D[1..m]$  of increasing positive
//      integers indicating the coin denominations where  $D[1] = 1$ 
//Output: The minimum number of coins that add up to  $n$ 
 $F[0] \leftarrow 0$ 
for  $i \leftarrow 1$  to  $n$  do
     $temp \leftarrow \infty$ ;  $j \leftarrow 1$ 
    while  $j \leq m$  and  $i \geq D[j]$  do
         $temp \leftarrow \min(F[i - D[j]], temp)$ 
         $j \leftarrow j + 1$ 
     $F[i] \leftarrow temp + 1$ 
return  $F[n]$ 
```

The application of the algorithm to amount $n = 6$ and denominations 1, 3, 4 is shown in Figure 8.2. The answer it yields is two coins. The time and space efficiencies of the algorithm are obviously $O(nm)$ and $\Theta(n)$, respectively.

above the cells in the first row, and there are no adjacent cells to the left of the cells in the first column. For those cells, we assume that $F(i-1, j)$ and $F(i, j-1)$ are equal to 0 for their nonexistent neighbors. Therefore, the largest number of coins the robot can bring to cell (i, j) is the maximum of these two numbers plus one possible coin at cell (i, j) itself. In other words, we have the following formula for $F(i, j)$:

$$\begin{aligned} F(i, j) &= \max\{F(i-1, j), F(i, j-1)\} + c_{ij} \quad \text{for } 1 \leq i \leq n, \quad 1 \leq j \leq m \\ F(0, j) &= 0 \quad \text{for } 1 \leq j \leq m \quad \text{and} \quad F(i, 0) = 0 \quad \text{for } 1 \leq i \leq n, \end{aligned} \quad (8.5)$$

where $c_{ij} = 1$ if there is a coin in cell (i, j) , and $c_{ij} = 0$ otherwise.

Using these formulas, we can fill in the $n \times m$ table of $F(i, j)$ values either row by row or column by column, as is typical for dynamic programming algorithms involving two-dimensional tables.

ALGORITHM *RobotCoinCollection*($C[1..n, 1..m]$)

```
//Applies dynamic programming to compute the largest number of
//coins a robot can collect on an  $n \times m$  board by starting at (1, 1)
//and moving right and down from upper left to down right corner
//Input: Matrix  $C[1..n, 1..m]$  whose elements are equal to 1 and 0
//for cells with and without a coin, respectively
//Output: Largest number of coins the robot can bring to cell  $(n, m)$ 
 $F[1, 1] \leftarrow C[1, 1]$ ; for  $j \leftarrow 2$  to  $m$  do  $F[1, j] \leftarrow F[1, j-1] + C[1, j]$ 
for  $i \leftarrow 2$  to  $n$  do
     $F[i, 1] \leftarrow F[i-1, 1] + C[i, 1]$ 
    for  $j \leftarrow 2$  to  $m$  do
         $F[i, j] \leftarrow \max(F[i-1, j], F[i, j-1]) + C[i, j]$ 
return  $F[n, m]$ 
```

The algorithm is illustrated in Figure 8.3b for the coin setup in Figure 8.3a. Since computing the value of $F(i, j)$ by formula (8.5) for each cell of the table takes constant time, the time efficiency of the algorithm is $\Theta(nm)$. Its space efficiency is, obviously, also $\Theta(nm)$.

Tracing the computations backward makes it possible to get an optimal path: if $F(i-1, j) > F(i, j-1)$, an optimal path to cell (i, j) must come down from the adjacent cell above it; if $F(i-1, j) < F(i, j-1)$, an optimal path to cell (i, j) must come from the adjacent cell on the left; and if $F(i-1, j) = F(i, j-1)$, it can reach cell (i, j) from either direction. This yields two optimal paths for the instance in Figure 8.3a, which are shown in Figure 8.3c. If ties are ignored, one optimal path can be obtained in $\Theta(n+m)$ time.

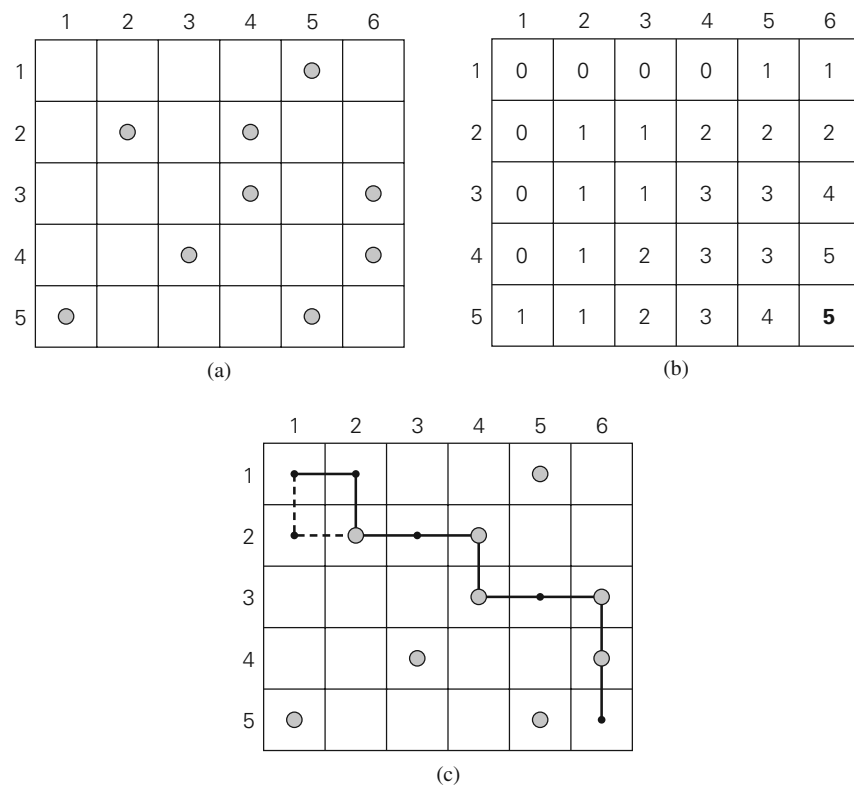


FIGURE 8.3 (a) Coins to collect. (b) Dynamic programming algorithm results. (c) Two paths to collect 5 coins, the maximum number of coins possible.

Exercises 8.1

1. What does dynamic programming have in common with divide-and-conquer? What is a principal difference between them?
2. Solve the instance 5, 1, 2, 10, 6 of the coin-row problem.
3. **a.** Show that the time efficiency of solving the coin-row problem by straight-forward application of recurrence (8.3) is exponential.
b. Show that the time efficiency of solving the coin-row problem by exhaustive search is at least exponential.
4. Apply the dynamic programming algorithm to find all the solutions to the change-making problem for the denominations 1, 3, 5 and the amount $n = 9$.

5. How would you modify the dynamic programming algorithm for the coin-collecting problem if some cells on the board are inaccessible for the robot? Apply your algorithm to the board below, where the inaccessible cells are shown by X's. How many optimal paths are there for this board?

	1	2	3	4	5	6
1		X		●		
2	●			X	●	
3		●		X	●	
4				●		●
5	X	X	X		●	

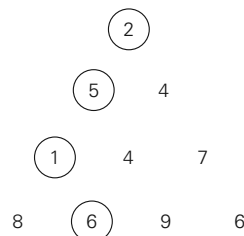
6. *Rod-cutting problem* Design a dynamic programming algorithm for the following problem. Find the maximum total sale price that can be obtained by cutting a rod of n units long into integer-length pieces if the sale price of a piece i units long is p_i for $i = 1, 2, \dots, n$. What are the time and space efficiencies of your algorithm?



7. *Shortest-path counting* A chess rook can move horizontally or vertically to any square in the same row or in the same column of a chessboard. Find the number of shortest paths by which a rook can move from one corner of a chessboard to the diagonally opposite corner. The length of a path is measured by the number of squares it passes through, including the first and the last squares. Solve the problem

- by a dynamic programming algorithm.
- by using elementary combinatorics.

8. *Minimum-sum descent* Some positive integers are arranged in an equilateral triangle with n numbers in its base like the one shown in the figure below for $n = 4$. The problem is to find the smallest sum in a descent from the triangle apex to its base through a sequence of adjacent numbers (shown in the figure by the circles). Design a dynamic programming algorithm for this problem and indicate its time efficiency.



9. *Binomial coefficient* Design an efficient algorithm for computing the binomial coefficient $C(n, k)$ that uses no multiplications. What are the time and space efficiencies of your algorithm?
10. *Longest path in a dag*
 - a. Design an efficient algorithm for finding the length of the longest path in a dag. (This problem is important both as a prototype of many other dynamic programming applications and in its own right because it determines the minimal time needed for completing a project comprising precedence-constrained tasks.)
 - b. Show how to reduce the coin-row problem discussed in this section to the problem of finding a longest path in a dag.
11. *Maximum square submatrix* Given an $m \times n$ boolean matrix B , find its largest square submatrix whose elements are all zeros. Design a dynamic programming algorithm and indicate its time efficiency. (The algorithm may be useful for, say, finding the largest free square area on a computer screen or for selecting a construction site.)
12. *World Series odds* Consider two teams, A and B , playing a series of games until one of the teams wins n games. Assume that the probability of A winning a game is the same for each game and equal to p , and the probability of A losing a game is $q = 1 - p$. (Hence, there are no ties.) Let $P(i, j)$ be the probability of A winning the series if A needs i more games to win the series and B needs j more games to win the series.
 - a. Set up a recurrence relation for $P(i, j)$ that can be used by a dynamic programming algorithm.
 - b. Find the probability of team A winning a seven-game series if the probability of it winning a game is 0.4.
 - c. Write pseudocode of the dynamic programming algorithm for solving this problem and determine its time and space efficiencies.

8.2 The Knapsack Problem and Memory Functions

We start this section with designing a dynamic programming algorithm for the knapsack problem: given n items of known weights w_1, \dots, w_n and values v_1, \dots, v_n and a knapsack of capacity W , find the most valuable subset of the items that fit into the knapsack. (This problem was introduced in Section 3.4, where we discussed solving it by exhaustive search.) We assume here that all the weights and the knapsack capacity are positive integers; the item values do not have to be integers.

To design a dynamic programming algorithm, we need to derive a recurrence relation that expresses a solution to an instance of the knapsack problem in terms

of solutions to its smaller subinstances. Let us consider an instance defined by the first i items, $1 \leq i \leq n$, with weights w_1, \dots, w_i , values v_1, \dots, v_i , and knapsack capacity j , $1 \leq j \leq W$. Let $F(i, j)$ be the value of an optimal solution to this instance, i.e., the value of the most valuable subset of the first i items that fit into the knapsack of capacity j . We can divide all the subsets of the first i items that fit the knapsack of capacity j into two categories: those that do not include the i th item and those that do. Note the following:

1. Among the subsets that do not include the i th item, the value of an optimal subset is, by definition, $F(i-1, j)$.
2. Among the subsets that do include the i th item (hence, $j - w_i \geq 0$), an optimal subset is made up of this item and an optimal subset of the first $i-1$ items that fits into the knapsack of capacity $j - w_i$. The value of such an optimal subset is $v_i + F(i-1, j - w_i)$.

Thus, the value of an optimal solution among all feasible subsets of the first i items is the maximum of these two values. Of course, if the i th item does not fit into the knapsack, the value of an optimal subset selected from the first i items is the same as the value of an optimal subset selected from the first $i-1$ items. These observations lead to the following recurrence:

$$F(i, j) = \begin{cases} \max\{F(i-1, j), v_i + F(i-1, j - w_i)\} & \text{if } j - w_i \geq 0, \\ F(i-1, j) & \text{if } j - w_i < 0. \end{cases} \quad (8.6)$$

It is convenient to define the initial conditions as follows:

$$F(0, j) = 0 \text{ for } j \geq 0 \quad \text{and} \quad F(i, 0) = 0 \text{ for } i \geq 0. \quad (8.7)$$

Our goal is to find $F(n, W)$, the maximal value of a subset of the n given items that fit into the knapsack of capacity W , and an optimal subset itself.

Figure 8.4 illustrates the values involved in equations (8.6) and (8.7). For $i, j > 0$, to compute the entry in the i th row and the j th column, $F(i, j)$, we compute the maximum of the entry in the previous row and the same column and the sum of v_i and the entry in the previous row and w_i columns to the left. The table can be filled either row by row or column by column.

	0	$j - w_i$	j	W
0	0	0	0	0
$i-1$	0	$F(i-1, j - w_i)$	$F(i-1, j)$	
w_i, v_i i	0		$F(i, j)$	
n	0			goal

FIGURE 8.4 Table for solving the knapsack problem by dynamic programming.

		capacity j					
	i	0	1	2	3	4	5
	0	0	0	0	0	0	0
$w_1 = 2, v_1 = 12$	1	0	0	12	12	12	12
$w_2 = 1, v_2 = 10$	2	0	10	12	22	22	22
$w_3 = 3, v_3 = 20$	3	0	10	12	22	30	32
$w_4 = 2, v_4 = 15$	4	0	10	15	25	30	37

FIGURE 8.5 Example of solving an instance of the knapsack problem by the dynamic programming algorithm.

EXAMPLE 1 Let us consider the instance given by the following data:

item	weight	value	capacity $W = 5$.
1	2	\$12	
2	1	\$10	
3	3	\$20	
4	2	\$15	

The dynamic programming table, filled by applying formulas (8.6) and (8.7), is shown in Figure 8.5.

Thus, the maximal value is $F(4, 5) = \$37$. We can find the composition of an optimal subset by backtracing the computations of this entry in the table. Since $F(4, 5) > F(3, 5)$, item 4 has to be included in an optimal solution along with an optimal subset for filling $5 - 2 = 3$ remaining units of the knapsack capacity. The value of the latter is $F(3, 3)$. Since $F(3, 3) = F(2, 3)$, item 3 need not be in an optimal subset. Since $F(2, 3) > F(1, 3)$, item 2 is a part of an optimal selection, which leaves element $F(1, 3 - 1)$ to specify its remaining composition. Similarly, since $F(1, 2) > F(0, 2)$, item 1 is the final part of the optimal solution {item 1, item 2, item 4}. ■

The time efficiency and space efficiency of this algorithm are both in $\Theta(nW)$. The time needed to find the composition of an optimal solution is in $O(n)$. You are asked to prove these assertions in the exercises.

Memory Functions

As we discussed at the beginning of this chapter and illustrated in subsequent sections, dynamic programming deals with problems whose solutions satisfy a recurrence relation with overlapping subproblems. The direct top-down approach to finding a solution to such a recurrence leads to an algorithm that solves common subproblems more than once and hence is very inefficient (typically, exponential

or worse). The classic dynamic programming approach, on the other hand, works bottom up: it fills a table with solutions to *all* smaller subproblems, but each of them is solved only once. An unsatisfying aspect of this approach is that solutions to some of these smaller subproblems are often not necessary for getting a solution to the problem given. Since this drawback is not present in the top-down approach, it is natural to try to combine the strengths of the top-down and bottom-up approaches. The goal is to get a method that solves only subproblems that are necessary and does so only once. Such a method exists; it is based on using **memory functions**.

This method solves a given problem in the top-down manner but, in addition, maintains a table of the kind that would have been used by a bottom-up dynamic programming algorithm. Initially, all the table's entries are initialized with a special "null" symbol to indicate that they have not yet been calculated. Thereafter, whenever a new value needs to be calculated, the method checks the corresponding entry in the table first: if this entry is not "null," it is simply retrieved from the table; otherwise, it is computed by the recursive call whose result is then recorded in the table.

The following algorithm implements this idea for the knapsack problem. After initializing the table, the recursive function needs to be called with $i = n$ (the number of items) and $j = W$ (the knapsack capacity).

ALGORITHM *MFKnapsack*(i, j)

```
//Implements the memory function method for the knapsack problem
//Input: A nonnegative integer  $i$  indicating the number of the first
//       items being considered and a nonnegative integer  $j$  indicating
//       the knapsack capacity
//Output: The value of an optimal feasible subset of the first  $i$  items
//Note: Uses as global variables input arrays  $Weights[1..n]$ ,  $Values[1..n]$ ,
//and table  $F[0..n, 0..W]$  whose entries are initialized with  $-1$ 's except for
//row 0 and column 0 initialized with 0's
if  $F[i, j] < 0$ 
    if  $j < Weights[i]$ 
         $value \leftarrow MFKnapsack(i - 1, j)$ 
    else
         $value \leftarrow \max(MFKnapsack(i - 1, j),$ 
                           $Values[i] + MFKnapsack(i - 1, j - Weights[i]))$ 
     $F[i, j] \leftarrow value$ 
return  $F[i, j]$ 
```

EXAMPLE 2 Let us apply the memory function method to the instance considered in Example 1. The table in Figure 8.6 gives the results. Only 11 out of 20 nontrivial values (i.e., not those in row 0 or in column 0) have been computed.

		capacity j						
		i	0	1	2	3	4	5
		0	0	0	0	0	0	0
$w_1 = 2, v_1 = 12$	1	0	0	12	12	12	12	
$w_2 = 1, v_2 = 10$	2	0	—	12	22	—	22	
$w_3 = 3, v_3 = 20$	3	0	—	—	22	—	32	
$w_4 = 2, v_4 = 15$	4	0	—	—	—	—	37	

FIGURE 8.6 Example of solving an instance of the knapsack problem by the memory function algorithm.

Just one nontrivial entry, $V(1, 2)$, is retrieved rather than being recomputed. For larger instances, the proportion of such entries can be significantly larger. ■

In general, we cannot expect more than a constant-factor gain in using the memory function method for the knapsack problem, because its time efficiency class is the same as that of the bottom-up algorithm (why?). A more significant improvement can be expected for dynamic programming algorithms in which a computation of one value takes more than constant time. You should also keep in mind that a memory function algorithm may be less space-efficient than a space-efficient version of a bottom-up algorithm.

Exercises 8.2

1. a. Apply the bottom-up dynamic programming algorithm to the following instance of the knapsack problem:

item	weight	value	capacity $W = 6$.
1	3	\$25	
2	2	\$20	
3	1	\$15	
4	4	\$40	
5	5	\$50	

- b. How many different optimal subsets does the instance of part (a) have?
- c. In general, how can we use the table generated by the dynamic programming algorithm to tell whether there is more than one optimal subset for the knapsack problem's instance?

2.
 - a. Write pseudocode of the bottom-up dynamic programming algorithm for the knapsack problem.
 - b. Write pseudocode of the algorithm that finds the composition of an optimal subset from the table generated by the bottom-up dynamic programming algorithm for the knapsack problem.
3. For the bottom-up dynamic programming algorithm for the knapsack problem, prove that
 - a. its time efficiency is $\Theta(nW)$.
 - b. its space efficiency is $\Theta(nW)$.
 - c. the time needed to find the composition of an optimal subset from a filled dynamic programming table is $O(n)$.
4.
 - a. True or false: A sequence of values in a row of the dynamic programming table for the knapsack problem is always nondecreasing?
 - b. True or false: A sequence of values in a column of the dynamic programming table for the knapsack problem is always nondecreasing?
5. Design a dynamic programming algorithm for the version of the knapsack problem in which there are unlimited quantities of copies for each of the n item kinds given. Indicate the time efficiency of the algorithm.
6. Apply the memory function method to the instance of the knapsack problem given in Problem 1. Indicate the entries of the dynamic programming table that are (i) never computed by the memory function method, (ii) retrieved without a recomputation.
7. Prove that the efficiency class of the memory function algorithm for the knapsack problem is the same as that of the bottom-up algorithm (see Problem 3).
8. Explain why the memory function approach is unattractive for the problem of computing a binomial coefficient by the formula $C(n, k) = C(n-1, k-1) + C(n-1, k)$.
9. Write a research report on one of the following well-known applications of dynamic programming:
 - a. finding the longest common subsequence in two sequences
 - b. optimal string editing
 - c. minimal triangulation of a polygon

8.3 Optimal Binary Search Trees

A binary search tree is one of the most important data structures in computer science. One of its principal applications is to implement a dictionary, a set of elements with the operations of searching, insertion, and deletion. If probabilities

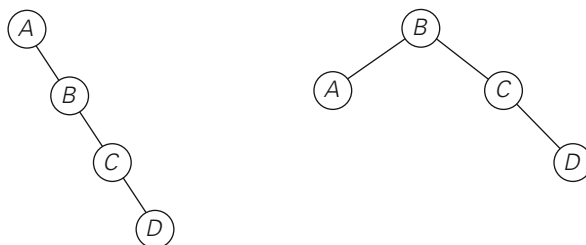


FIGURE 8.7 Two out of 14 possible binary search trees with keys A , B , C , and D .

of searching for elements of a set are known—e.g., from accumulated data about past searches—it is natural to pose a question about an optimal binary search tree for which the average number of comparisons in a search is the smallest possible. For simplicity, we limit our discussion to minimizing the average number of comparisons in a successful search. The method can be extended to include unsuccessful searches as well.

As an example, consider four keys A , B , C , and D to be searched for with probabilities 0.1, 0.2, 0.4, and 0.3, respectively. Figure 8.7 depicts two out of 14 possible binary search trees containing these keys. The average number of comparisons in a successful search in the first of these trees is $0.1 \cdot 1 + 0.2 \cdot 2 + 0.4 \cdot 3 + 0.3 \cdot 4 = 2.9$, and for the second one it is $0.1 \cdot 2 + 0.2 \cdot 1 + 0.4 \cdot 2 + 0.3 \cdot 3 = 2.1$. Neither of these two trees is, in fact, optimal. (Can you tell which binary tree is optimal?)

For our tiny example, we could find the optimal tree by generating all 14 binary search trees with these keys. As a general algorithm, this exhaustive-search approach is unrealistic: the total number of binary search trees with n keys is equal to the n th **Catalan number**,

$$c(n) = \frac{1}{n+1} \binom{2n}{n} \quad \text{for } n > 0, \quad c(0) = 1,$$

which grows to infinity as fast as $4^n/n^{1.5}$ (see Problem 7 in this section's exercises).

So let a_1, \dots, a_n be distinct keys ordered from the smallest to the largest and let p_1, \dots, p_n be the probabilities of searching for them. Let $C(i, j)$ be the smallest average number of comparisons made in a successful search in a binary search tree T_i^j made up of keys a_i, \dots, a_j , where i, j are some integer indices, $1 \leq i \leq j \leq n$. Following the classic dynamic programming approach, we will find values of $C(i, j)$ for all smaller instances of the problem, although we are interested just in $C(1, n)$. To derive a recurrence underlying a dynamic programming algorithm, we will consider all possible ways to choose a root a_k among the keys a_i, \dots, a_j . For such a binary search tree (Figure 8.8), the root contains key a_k , the left subtree T_i^{k-1} contains keys a_i, \dots, a_{k-1} optimally arranged, and the right subtree T_{k+1}^j

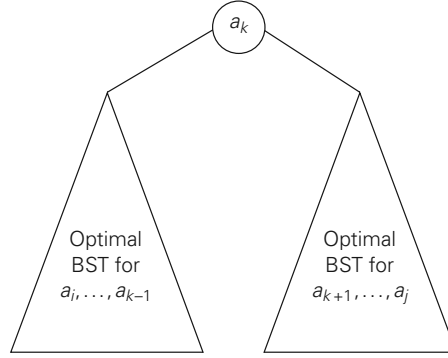


FIGURE 8.8 Binary search tree (BST) with root a_k and two optimal binary search subtrees T_i^{k-1} and T_{k+1}^j .

contains keys a_{k+1}, \dots, a_j also optimally arranged. (Note how we are taking advantage of the principle of optimality here.)

If we count tree levels starting with 1 to make the comparison numbers equal the keys' levels, the following recurrence relation is obtained:

$$\begin{aligned}
 C(i, j) &= \min_{i \leq k \leq j} \left\{ p_k \cdot 1 + \sum_{s=i}^{k-1} p_s \cdot (\text{level of } a_s \text{ in } T_i^{k-1} + 1) \right. \\
 &\quad \left. + \sum_{s=k+1}^j p_s \cdot (\text{level of } a_s \text{ in } T_{k+1}^j + 1) \right\} \\
 &= \min_{i \leq k \leq j} \left\{ \sum_{s=i}^{k-1} p_s \cdot \text{level of } a_s \text{ in } T_i^{k-1} + \sum_{s=k+1}^j p_s \cdot \text{level of } a_s \text{ in } T_{k+1}^j + \sum_{s=i}^j p_s \right\} \\
 &= \min_{i \leq k \leq j} \{ C(i, k-1) + C(k+1, j) \} + \sum_{s=i}^j p_s.
 \end{aligned}$$

Thus, we have the recurrence

$$C(i, j) = \min_{i \leq k \leq j} \{ C(i, k-1) + C(k+1, j) \} + \sum_{s=i}^j p_s \quad \text{for } 1 \leq i \leq j \leq n. \quad (8.8)$$

We assume in formula (8.8) that $C(i, i-1) = 0$ for $1 \leq i \leq n+1$, which can be interpreted as the number of comparisons in the empty tree. Note that this formula implies that

$$C(i, i) = p_i \quad \text{for } 1 \leq i \leq n,$$

as it should be for a one-node binary search tree containing a_i .

	0	1				j	n
1	0	p_1					goal
		0	p_2				
i						$C[i, j]$	
							p_n
$n+1$							0

FIGURE 8.9 Table of the dynamic programming algorithm for constructing an optimal binary search tree.

The two-dimensional table in Figure 8.9 shows the values needed for computing $C(i, j)$ by formula (8.8): they are in row i and the columns to the left of column j and in column j and the rows below row i . The arrows point to the pairs of entries whose sums are computed in order to find the smallest one to be recorded as the value of $C(i, j)$. This suggests filling the table along its diagonals, starting with all zeros on the main diagonal and given probabilities p_i , $1 \leq i \leq n$, right above it and moving toward the upper right corner.

The algorithm we just sketched computes $C(1, n)$ —the average number of comparisons for successful searches in the optimal binary tree. If we also want to get the optimal tree itself, we need to maintain another two-dimensional table to record the value of k for which the minimum in (8.8) is achieved. The table has the same shape as the table in Figure 8.9 and is filled in the same manner, starting with entries $R(i, i) = i$ for $1 \leq i \leq n$. When the table is filled, its entries indicate indices of the roots of the optimal subtrees, which makes it possible to reconstruct an optimal tree for the entire set given.

EXAMPLE Let us illustrate the algorithm by applying it to the four-key set we used at the beginning of this section:

key	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>
probability	0.1	0.2	0.4	0.3

The initial tables look like this:

main table						root table					
	0	1	2	3	4		0	1	2	3	4
1	0	0.1				1		1			
2		0	0.2			2			2		
3			0	0.4		3				3	
4				0	0.3	4					4
5					0	5					

Let us compute $C(1, 2)$:

$$C(1, 2) = \min \left\{ \begin{array}{l} k=1: C(1, 0) + C(2, 2) + \sum_{s=1}^2 p_s = 0 + 0.2 + 0.3 = 0.5 \\ k=2: C(1, 1) + C(3, 2) + \sum_{s=1}^2 p_s = 0.1 + 0 + 0.3 = 0.4 \end{array} \right\} = 0.4.$$

Thus, out of two possible binary trees containing the first two keys, *A* and *B*, the root of the optimal tree has index 2 (i.e., it contains *B*), and the average number of comparisons in a successful search in this tree is 0.4.

We will ask you to finish the computations in the exercises. You should arrive at the following final tables:

main table						root table					
	0	1	2	3	4		0	1	2	3	4
1	0	0.1	0.4	1.1	1.7	1		1	2	3	3
2		0	0.2	0.8	1.4	2			2	3	3
3			0	0.4	1.0	3				3	3
4				0	0.3	4					4
5					0	5					

Thus, the average number of key comparisons in the optimal tree is equal to 1.7. Since $R(1, 4) = 3$, the root of the optimal tree contains the third key, i.e., *C*. Its left subtree is made up of keys *A* and *B*, and its right subtree contains just key *D* (why?). To find the specific structure of these subtrees, we find first their roots by consulting the root table again as follows. Since $R(1, 2) = 2$, the root of the optimal tree containing *A* and *B* is *B*, with *A* being its left child (and the root of the one-node tree: $R(1, 1) = 1$). Since $R(4, 4) = 4$, the root of this one-node optimal tree is its only key *D*. Figure 8.10 presents the optimal tree in its entirety.

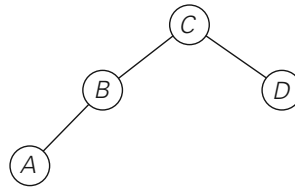


FIGURE 8.10 Optimal binary search tree for the example.

Here is pseudocode of the dynamic programming algorithm.

ALGORITHM *OptimalBST*($P[1..n]$)

```

//Finds an optimal binary search tree by dynamic programming
//Input: An array  $P[1..n]$  of search probabilities for a sorted list of  $n$  keys
//Output: Average number of comparisons in successful searches in the
//        optimal BST and table  $R$  of subtrees' roots in the optimal BST
for  $i \leftarrow 1$  to  $n$  do
     $C[i, i - 1] \leftarrow 0$ 
     $C[i, i] \leftarrow P[i]$ 
     $R[i, i] \leftarrow i$ 
 $C[n + 1, n] \leftarrow 0$ 
for  $d \leftarrow 1$  to  $n - 1$  do //diagonal count
    for  $i \leftarrow 1$  to  $n - d$  do
         $j \leftarrow i + d$ 
         $minval \leftarrow \infty$ 
        for  $k \leftarrow i$  to  $j$  do
            if  $C[i, k - 1] + C[k + 1, j] < minval$ 
                 $minval \leftarrow C[i, k - 1] + C[k + 1, j]; \quad kmin \leftarrow k$ 
             $R[i, j] \leftarrow kmin$ 
         $sum \leftarrow P[i]; \quad \textbf{for } s \leftarrow i + 1 \textbf{ to } j \textbf{ do } sum \leftarrow sum + P[s]$ 
         $C[i, j] \leftarrow minval + sum$ 
return  $C[1, n], R$ 
  
```

The algorithm's space efficiency is clearly quadratic; the time efficiency of this version of the algorithm is cubic (why?). A more careful analysis shows that entries in the root table are always nondecreasing along each row and column. This limits values for $R(i, j)$ to the range $R(i, j - 1), \dots, R(i + 1, j)$ and makes it possible to reduce the running time of the algorithm to $\Theta(n^2)$.

Exercises 8.3

1. Finish the computations started in the section's example of constructing an optimal binary search tree.
2. **a.** Why is the time efficiency of algorithm *OptimalBST* cubic?
b. Why is the space efficiency of algorithm *OptimalBST* quadratic?
3. Write pseudocode for a linear-time algorithm that generates the optimal binary search tree from the root table.
4. Devise a way to compute the sums $\sum_{s=i}^j p_s$, which are used in the dynamic programming algorithm for constructing an optimal binary search tree, in constant time (per sum).
5. True or false: The root of an optimal binary search tree always contains the key with the highest search probability?
6. How would you construct an optimal binary search tree for a set of n keys if all the keys are equally likely to be searched for? What will be the average number of comparisons in a successful search in such a tree if $n = 2^k$?
7. **a.** Show that the number of distinct binary search trees $b(n)$ that can be constructed for a set of n orderable keys satisfies the recurrence relation

$$b(n) = \sum_{k=0}^{n-1} b(k)b(n-1-k) \quad \text{for } n > 0, \quad b(0) = 1.$$

- b.** It is known that the solution to this recurrence is given by the Catalan numbers. Verify this assertion for $n = 1, 2, \dots, 5$.
- c.** Find the order of growth of $b(n)$. What implication does the answer to this question have for the exhaustive-search algorithm for constructing an optimal binary search tree?
8. Design a $\Theta(n^2)$ algorithm for finding an optimal binary search tree.
9. Generalize the optimal binary search algorithm by taking into account unsuccessful searches.
10. Write pseudocode of a memory function for the optimal binary search tree problem. You may limit your function to finding the smallest number of key comparisons in a successful search.
11. **Matrix chain multiplication** Consider the problem of minimizing the total number of multiplications made in computing the product of n matrices

$$A_1 \cdot A_2 \cdot \dots \cdot A_n$$

whose dimensions are $d_0 \times d_1, d_1 \times d_2, \dots, d_{n-1} \times d_n$, respectively. Assume that all intermediate products of two matrices are computed by the brute-force (definition-based) algorithm.

- a. Give an example of three matrices for which the number of multiplications in $(A_1 \cdot A_2) \cdot A_3$ and $A_1 \cdot (A_2 \cdot A_3)$ differ at least by a factor of 1000.
- b. How many different ways are there to compute the product of n matrices?
- c. Design a dynamic programming algorithm for finding an optimal order of multiplying n matrices.

8.4 Marshall's and Floyd's Algorithms

In this section, we look at two well-known algorithms: Marshall's algorithm for computing the transitive closure of a directed graph and Floyd's algorithm for the all-pairs shortest-paths problem. These algorithms are based on essentially the same idea: exploit a relationship between a problem and its simpler rather than smaller version. Marshall and Floyd published their algorithms without mentioning dynamic programming. Nevertheless, the algorithms certainly have a dynamic programming flavor and have come to be considered applications of this technique.

Marshall's Algorithm

Recall that the adjacency matrix $A = \{a_{ij}\}$ of a directed graph is the boolean matrix that has 1 in its i th row and j th column if and only if there is a directed edge from the i th vertex to the j th vertex. We may also be interested in a matrix containing the information about the existence of directed paths of arbitrary lengths between vertices of a given graph. Such a matrix, called the transitive closure of the digraph, would allow us to determine in constant time whether the j th vertex is reachable from the i th vertex.

Here are a few application examples. When a value in a spreadsheet cell is changed, the spreadsheet software must know all the other cells affected by the change. If the spreadsheet is modeled by a digraph whose vertices represent the spreadsheet cells and edges indicate cell dependencies, the transitive closure will provide such information. In software engineering, transitive closure can be used for investigating data flow and control flow dependencies as well as for inheritance testing of object-oriented software. In electronic engineering, it is used for redundancy identification and test generation for digital circuits.

DEFINITION The *transitive closure* of a directed graph with n vertices can be defined as the $n \times n$ boolean matrix $T = \{t_{ij}\}$, in which the element in the i th row and the j th column is 1 if there exists a nontrivial path (i.e., directed path of a positive length) from the i th vertex to the j th vertex; otherwise, t_{ij} is 0.

An example of a digraph, its adjacency matrix, and its transitive closure is given in Figure 8.11.

We can generate the transitive closure of a digraph with the help of depth-first search or breadth-first search. Performing either traversal starting at the i th

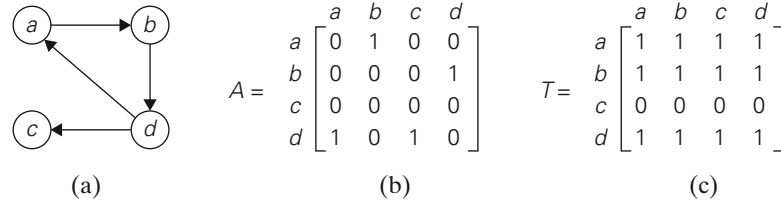


FIGURE 8.11 (a) Digraph. (b) Its adjacency matrix. (c) Its transitive closure.

vertex gives the information about the vertices reachable from it and hence the columns that contain 1's in the i th row of the transitive closure. Thus, doing such a traversal for every vertex as a starting point yields the transitive closure in its entirety.

Since this method traverses the same digraph several times, we should hope that a better algorithm can be found. Indeed, such an algorithm exists. It is called **Warshall's algorithm** after Stephen Warshall, who discovered it [War62]. It is convenient to assume that the digraph's vertices and hence the rows and columns of the adjacency matrix are numbered from 1 to n . Warshall's algorithm constructs the transitive closure through a series of $n \times n$ boolean matrices:

$$R^{(0)}, \dots, R^{(k-1)}, R^{(k)}, \dots, R^{(n)}. \quad (8.9)$$

Each of these matrices provides certain information about directed paths in the digraph. Specifically, the element $r_{ij}^{(k)}$ in the i th row and j th column of matrix $R^{(k)}$ ($i, j = 1, 2, \dots, n, k = 0, 1, \dots, n$) is equal to 1 if and only if there exists a directed path of a positive length from the i th vertex to the j th vertex with each intermediate vertex, if any, numbered not higher than k . Thus, the series starts with $R^{(0)}$, which does not allow any intermediate vertices in its paths; hence, $R^{(0)}$ is nothing other than the adjacency matrix of the digraph. (Recall that the adjacency matrix contains the information about one-edge paths, i.e., paths with no intermediate vertices.) $R^{(1)}$ contains the information about paths that can use the first vertex as intermediate; thus, with more freedom, so to speak, it may contain more 1's than $R^{(0)}$. In general, each subsequent matrix in series (8.9) has one more vertex to use as intermediate for its paths than its predecessor and hence may, but does not have to, contain more 1's. The last matrix in the series, $R^{(n)}$, reflects paths that can use all n vertices of the digraph as intermediate and hence is nothing other than the digraph's transitive closure.

The central point of the algorithm is that we can compute all the elements of each matrix $R^{(k)}$ from its immediate predecessor $R^{(k-1)}$ in series (8.9). Let $r_{ij}^{(k)}$, the element in the i th row and j th column of matrix $R^{(k)}$, be equal to 1. This means that there exists a path from the i th vertex v_i to the j th vertex v_j with each intermediate vertex numbered not higher than k :

$$v_i, \text{ a list of intermediate vertices each numbered not higher than } k, v_j. \quad (8.10)$$

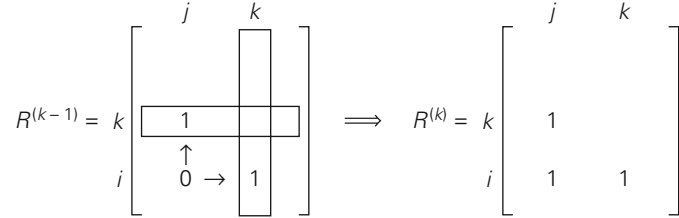


FIGURE 8.12 Rule for changing zeros in Warshall's algorithm.

Two situations regarding this path are possible. In the first, the list of its intermediate vertices does not contain the k th vertex. Then this path from v_i to v_j has intermediate vertices numbered not higher than $k - 1$, and therefore $r_{ij}^{(k-1)}$ is equal to 1 as well. The second possibility is that path (8.10) does contain the k th vertex v_k among the intermediate vertices. Without loss of generality, we may assume that v_k occurs only once in that list. (If it is not the case, we can create a new path from v_i to v_j with this property by simply eliminating all the vertices between the first and last occurrences of v_k in it.) With this caveat, path (8.10) can be rewritten as follows:

$$v_i, \text{ vertices numbered } \leq k - 1, v_k, \text{ vertices numbered } \leq k - 1, v_j.$$

The first part of this representation means that there exists a path from v_i to v_k with each intermediate vertex numbered not higher than $k - 1$ (hence, $r_{ik}^{(k-1)} = 1$), and the second part means that there exists a path from v_k to v_j with each intermediate vertex numbered not higher than $k - 1$ (hence, $r_{kj}^{(k-1)} = 1$).

What we have just proved is that if $r_{ij}^{(k)} = 1$, then either $r_{ij}^{(k-1)} = 1$ or both $r_{ik}^{(k-1)} = 1$ and $r_{kj}^{(k-1)} = 1$. It is easy to see that the converse of this assertion is also true. Thus, we have the following formula for generating the elements of matrix $R^{(k)}$ from the elements of matrix $R^{(k-1)}$:

$$r_{ij}^{(k)} = r_{ij}^{(k-1)} \quad \text{or} \quad \left(r_{ik}^{(k-1)} \text{ and } r_{kj}^{(k-1)} \right). \quad (8.11)$$

Formula (8.11) is at the heart of Warshall's algorithm. This formula implies the following rule for generating elements of matrix $R^{(k)}$ from elements of matrix $R^{(k-1)}$, which is particularly convenient for applying Warshall's algorithm by hand:

- If an element r_{ij} is 1 in $R^{(k-1)}$, it remains 1 in $R^{(k)}$.
- If an element r_{ij} is 0 in $R^{(k-1)}$, it has to be changed to 1 in $R^{(k)}$ if and only if the element in its row i and column k and the element in its column j and row k are both 1's in $R^{(k-1)}$. This rule is illustrated in Figure 8.12.

As an example, the application of Warshall's algorithm to the digraph in Figure 8.11 is shown in Figure 8.13.

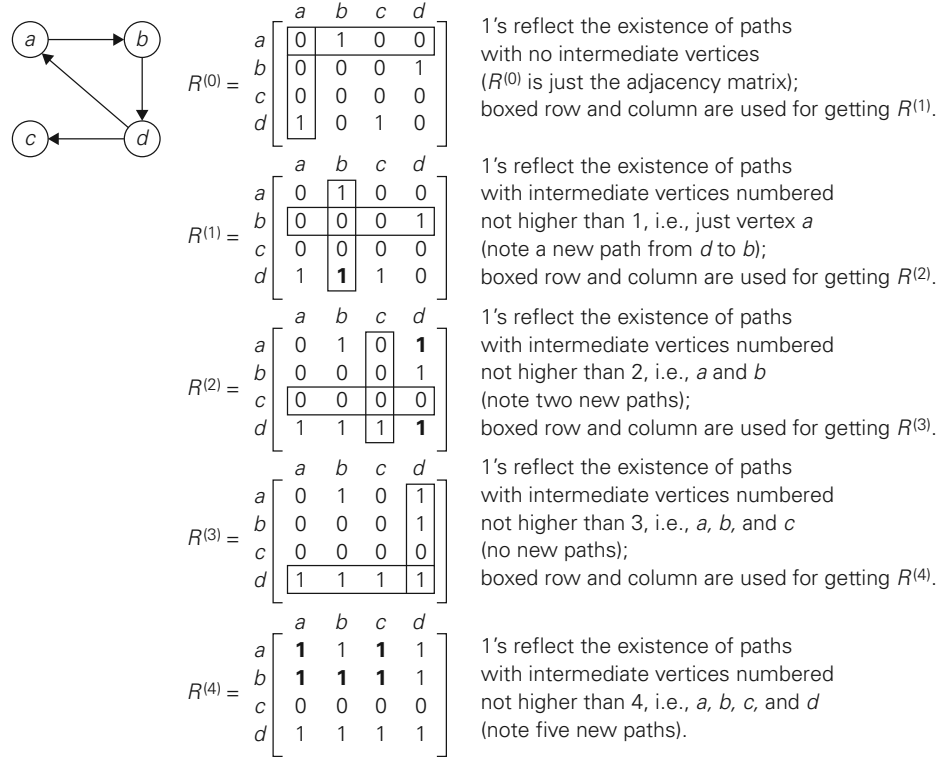


FIGURE 8.13 Application of Warshall's algorithm to the digraph shown. New 1's are in bold.

Here is pseudocode of Warshall's algorithm.

ALGORITHM *Warshall*($A[1..n, 1..n]$)

//Implements Warshall's algorithm for computing the transitive closure

//Input: The adjacency matrix A of a digraph with n vertices

//Output: The transitive closure of the digraph

$R^{(0)} \leftarrow A$

for $k \leftarrow 1$ **to** n **do**

for $i \leftarrow 1$ **to** n **do**

for $j \leftarrow 1$ **to** n **do**

$R^{(k)}[i, j] \leftarrow R^{(k-1)}[i, j] \text{ or } (R^{(k-1)}[i, k] \text{ and } R^{(k-1)}[k, j])$

return $R^{(n)}$

Several observations need to be made about Warshall's algorithm. First, it is remarkably succinct, is it not? Still, its time efficiency is only $\Theta(n^3)$. In fact, for sparse graphs represented by their adjacency lists, the traversal-based algorithm

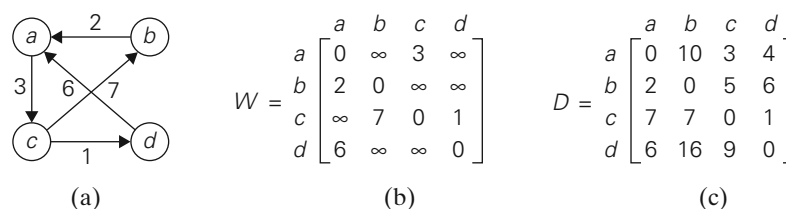


FIGURE 8.14 (a) Digraph. (b) Its weight matrix. (c) Its distance matrix.

mentioned at the beginning of this section has a better asymptotic efficiency than Warshall's algorithm (why?). We can speed up the above implementation of Warshall's algorithm for some inputs by restructuring its innermost loop (see Problem 4 in this section's exercises). Another way to make the algorithm run faster is to treat matrix rows as bit strings and employ the bitwise *or* operation available in most modern computer languages.

As to the space efficiency of Warshall's algorithm, the situation is similar to that of computing a Fibonacci number and some other dynamic programming algorithms. Although we used separate matrices for recording intermediate results of the algorithm, this is, in fact, unnecessary. Problem 3 in this section's exercises asks you to find a way of avoiding this wasteful use of the computer memory. Finally, we shall see below how the underlying idea of Warshall's algorithm can be applied to the more general problem of finding lengths of shortest paths in weighted graphs.

Floyd's Algorithm for the All-Pairs Shortest-Paths Problem

Given a weighted connected graph (undirected or directed), the ***all-pairs shortest-paths problem*** asks to find the distances—i.e., the lengths of the shortest paths—from each vertex to all other vertices. This is one of several variations of the problem involving shortest paths in graphs. Because of its important applications to communications, transportation networks, and operations research, it has been thoroughly studied over the years. Among recent applications of the all-pairs shortest-path problem is precomputing distances for motion planning in computer games.

It is convenient to record the lengths of shortest paths in an $n \times n$ matrix D called the ***distance matrix***: the element d_{ij} in the i th row and the j th column of this matrix indicates the length of the shortest path from the i th vertex to the j th vertex. For an example, see Figure 8.14.

We can generate the distance matrix with an algorithm that is very similar to Warshall's algorithm. It is called ***Floyd's algorithm*** after its co-inventor Robert W. Floyd.¹ It is applicable to both undirected and directed weighted graphs provided

1. Floyd explicitly referenced Warshall's paper in presenting his algorithm [Flo62]. Three years earlier, Bernard Roy published essentially the same algorithm in the proceedings of the French Academy of Sciences [Roy59].

that they do not contain a cycle of a negative length. (The distance between any two vertices in such a cycle can be made arbitrarily small by repeating the cycle enough times.) The algorithm can be enhanced to find not only the lengths of the shortest paths for all vertex pairs but also the shortest paths themselves (Problem 10 in this section's exercises).

Floyd's algorithm computes the distance matrix of a weighted graph with n vertices through a series of $n \times n$ matrices:

$$D^{(0)}, \dots, D^{(k-1)}, D^{(k)}, \dots, D^{(n)}. \quad (8.12)$$

Each of these matrices contains the lengths of shortest paths with certain constraints on the paths considered for the matrix in question. Specifically, the element $d_{ij}^{(k)}$ in the i th row and the j th column of matrix $D^{(k)}$ ($i, j = 1, 2, \dots, n$, $k = 0, 1, \dots, n$) is equal to the length of the shortest path among all paths from the i th vertex to the j th vertex with each intermediate vertex, if any, numbered not higher than k . In particular, the series starts with $D^{(0)}$, which does not allow any intermediate vertices in its paths; hence, $D^{(0)}$ is simply the weight matrix of the graph. The last matrix in the series, $D^{(n)}$, contains the lengths of the shortest paths among all paths that can use all n vertices as intermediate and hence is nothing other than the distance matrix being sought.

As in Warshall's algorithm, we can compute all the elements of each matrix $D^{(k)}$ from its immediate predecessor $D^{(k-1)}$ in series (8.12). Let $d_{ij}^{(k)}$ be the element in the i th row and the j th column of matrix $D^{(k)}$. This means that $d_{ij}^{(k)}$ is equal to the length of the shortest path among all paths from the i th vertex v_i to the j th vertex v_j with their intermediate vertices numbered not higher than k :

$$v_i, \text{ a list of intermediate vertices each numbered not higher than } k, v_j. \quad (8.13)$$

We can partition all such paths into two disjoint subsets: those that do not use the k th vertex v_k as intermediate and those that do. Since the paths of the first subset have their intermediate vertices numbered not higher than $k - 1$, the shortest of them is, by definition of our matrices, of length $d_{ij}^{(k-1)}$.

What is the length of the shortest path in the second subset? If the graph does not contain a cycle of a negative length, we can limit our attention only to the paths in the second subset that use vertex v_k as their intermediate vertex exactly once (because visiting v_k more than once can only increase the path's length). All such paths have the following form:

$$v_i, \text{ vertices numbered } \leq k - 1, v_k, \text{ vertices numbered } \leq k - 1, v_j.$$

In other words, each of the paths is made up of a path from v_i to v_k with each intermediate vertex numbered not higher than $k - 1$ and a path from v_k to v_j with each intermediate vertex numbered not higher than $k - 1$. The situation is depicted symbolically in Figure 8.15.

Since the length of the shortest path from v_i to v_k among the paths that use intermediate vertices numbered not higher than $k - 1$ is equal to $d_{ik}^{(k-1)}$ and the length of the shortest path from v_k to v_j among the paths that use intermediate

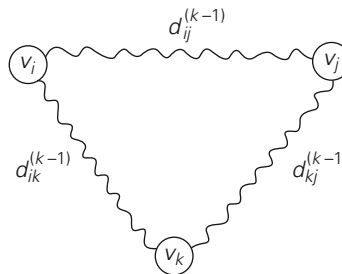


FIGURE 8.15 Underlying idea of Floyd's algorithm.

vertices numbered not higher than $k - 1$ is equal to $d_{kj}^{(k-1)}$, the length of the shortest path among the paths that use the k th vertex is equal to $d_{ik}^{(k-1)} + d_{kj}^{(k-1)}$. Taking into account the lengths of the shortest paths in both subsets leads to the following recurrence:

$$d_{ij}^{(k)} = \min\{d_{ij}^{(k-1)}, d_{ik}^{(k-1)} + d_{kj}^{(k-1)}\} \quad \text{for } k \geq 1, \quad d_{ij}^{(0)} = w_{ij}. \quad (8.14)$$

To put it another way, the element in row i and column j of the current distance matrix $D^{(k-1)}$ is replaced by the sum of the elements in the same row i and the column k and in the same column j and the row k if and only if the latter sum is smaller than its current value.

The application of Floyd's algorithm to the graph in Figure 8.14 is illustrated in Figure 8.16.

Here is pseudocode of Floyd's algorithm. It takes advantage of the fact that the next matrix in sequence (8.12) can be written over its predecessor.

ALGORITHM *Floyd*($W[1..n, 1..n]$)

//Implements Floyd's algorithm for the all-pairs shortest-paths problem

//Input: The weight matrix W of a graph with no negative-length cycle

//Output: The distance matrix of the shortest paths' lengths

$D \leftarrow W$ //is not necessary if W can be overwritten

for $k \leftarrow 1$ **to** n **do**

for $i \leftarrow 1$ **to** n **do**

for $j \leftarrow 1$ **to** n **do**

$D[i, j] \leftarrow \min\{D[i, j], D[i, k] + D[k, j]\}$

return D

Obviously, the time efficiency of Floyd's algorithm is cubic—as is the time efficiency of Warshall's algorithm. In the next chapter, we examine Dijkstra's algorithm—another method for finding shortest paths.

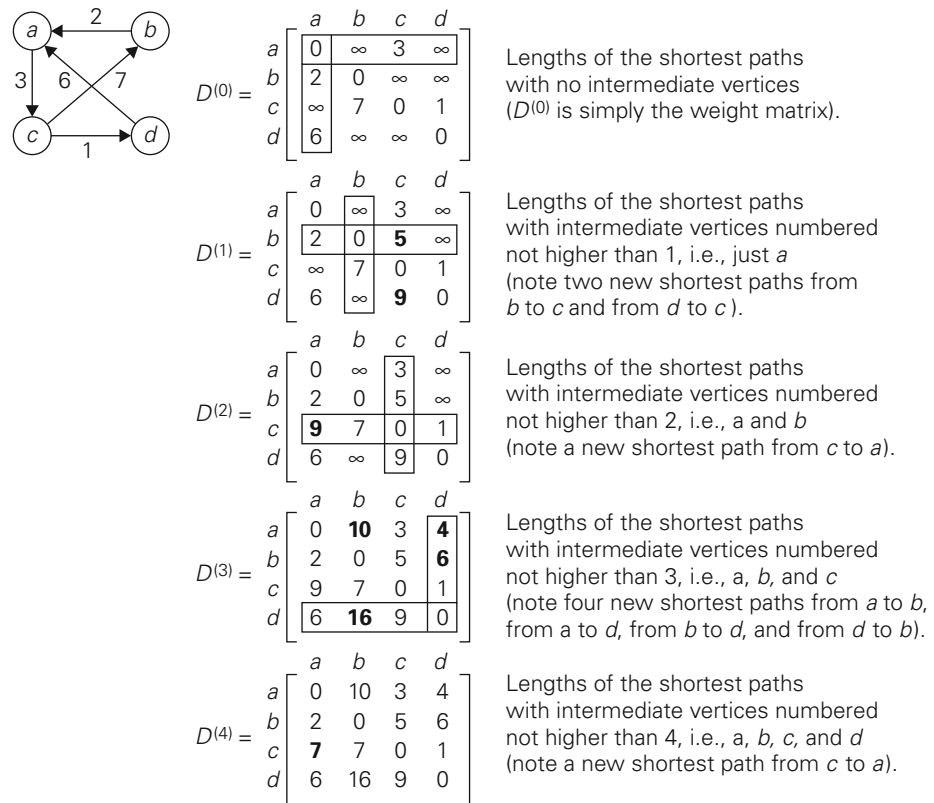


FIGURE 8.16 Application of Floyd's algorithm to the digraph shown. Updated elements are shown in bold.

Exercises 8.4


1. Apply Warshall's algorithm to find the transitive closure of the digraph defined by the following adjacency matrix:

$$\begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

2.
 - a. Prove that the time efficiency of Warshall's algorithm is cubic.
 - b. Explain why the time efficiency class of Warshall's algorithm is inferior to that of the traversal-based algorithm for sparse graphs represented by their adjacency lists.

3. Explain how to implement Warshall's algorithm without using extra memory for storing elements of the algorithm's intermediate matrices.
4. Explain how to restructure the innermost loop of the algorithm *Warshall* to make it run faster at least on some inputs.
5. Rewrite pseudocode of Warshall's algorithm assuming that the matrix rows are represented by bit strings on which the bitwise *or* operation can be performed.
6. a. Explain how Warshall's algorithm can be used to determine whether a given digraph is a dag (directed acyclic graph). Is it a good algorithm for this problem?
b. Is it a good idea to apply Warshall's algorithm to find the transitive closure of an undirected graph?
7. Solve the all-pairs shortest-path problem for the digraph with the following weight matrix:

$$\begin{bmatrix} 0 & 2 & \infty & 1 & 8 \\ 6 & 0 & 3 & 2 & \infty \\ \infty & \infty & 0 & 4 & \infty \\ \infty & \infty & 2 & 0 & 3 \\ 3 & \infty & \infty & \infty & 0 \end{bmatrix}$$

8. Prove that the next matrix in sequence (8.12) of Floyd's algorithm can be written over its predecessor.
9. Give an example of a graph or a digraph with negative weights for which Floyd's algorithm does not yield the correct result.
10. Enhance Floyd's algorithm so that shortest paths themselves, not just their lengths, can be found.
11.  *Jack Straws* In the game of Jack Straws, a number of plastic or wooden "straws" are dumped on the table and players try to remove them one by one without disturbing the other straws. Here, we are only concerned with whether various pairs of straws are connected by a path of touching straws. Given a list of the endpoints for $n > 1$ straws (as if they were dumped on a large piece of graph paper), determine all the pairs of straws that are connected. Note that touching is connecting, but also that two straws can be connected indirectly via other connected straws. [1994 East-Central Regionals of the ACM International Collegiate Programming Contest]

SUMMARY

- *Dynamic programming* is a technique for solving problems with overlapping subproblems. Typically, these subproblems arise from a recurrence relating a solution to a given problem with solutions to its smaller subproblems of the

same type. Dynamic programming suggests solving each smaller subproblem once and recording the results in a table from which a solution to the original problem can be then obtained.

- Applicability of dynamic programming to an optimization problem requires the problem to satisfy the *principle of optimality*: an optimal solution to any of its instances must be made up of optimal solutions to its subinstances.
- Among many other problems, the *change-making problem* with arbitrary coin denominations can be solved by dynamic programming.
- Solving a knapsack problem by a dynamic programming algorithm exemplifies an application of this technique to difficult problems of combinatorial optimization.
- The *memory function* technique seeks to combine the strengths of the top-down and bottom-up approaches to solving problems with overlapping subproblems. It does this by solving, in the top-down fashion but only once, just the necessary subproblems of a given problem and recording their solutions in a table.
- Dynamic programming can be used for constructing an *optimal binary search tree* for a given set of keys and known probabilities of searching for them.
- *Warshall's algorithm* for finding the *transitive closure* and *Floyd's algorithm* for the *all-pairs shortest-paths problem* are based on the idea that can be interpreted as an application of the dynamic programming technique.