

# MODULE-4

## 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:

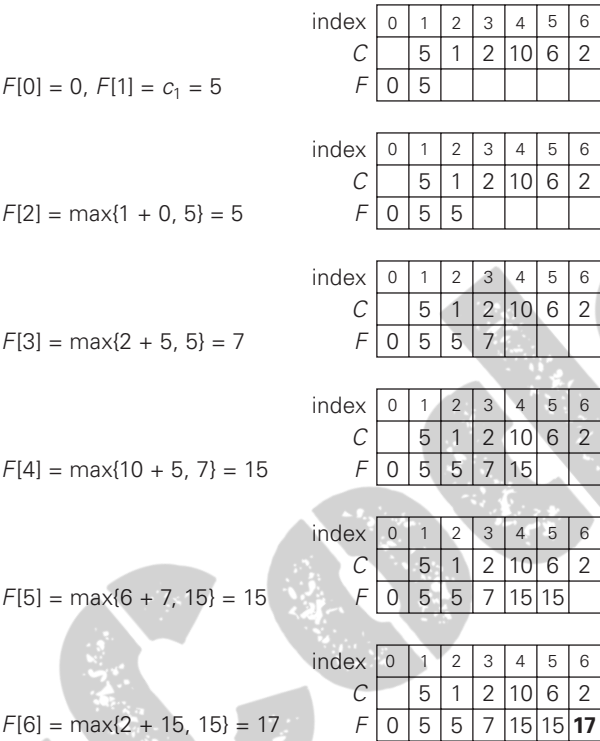
$$\begin{aligned} F(n) &= \max\{c_n + F(n - 2), F(n - 1)\} \quad \text{for } n > 1, \\ F(0) &= 0, \quad F(1) = c_1. \end{aligned} \tag{8.3}$$

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



**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  $F(n)$ :

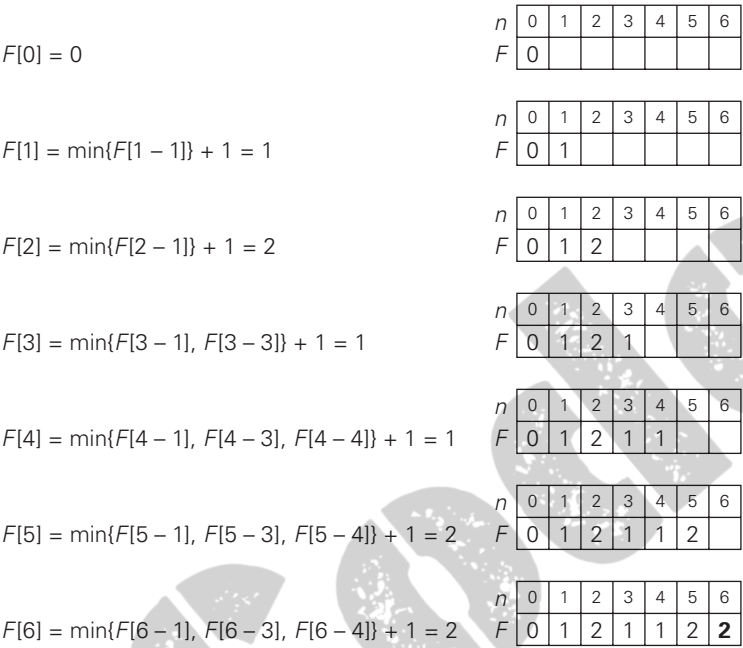
$$\begin{aligned} F(n) &= \min_{j: n \geq d_j} \{F(n - d_j)\} + 1 \quad \text{for } n > 0, \\ F(0) &= 0. \end{aligned} \tag{8.4}$$

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.



**FIGURE 8.2** Application of Algorithm *MinCoinChange* to amount  $n = 6$  and coin denominations 1, 3, and 4.

To find the coins of an optimal solution, we need to backtrack the computations to see which of the denominations produced the minima in formula (8.4). For the instance considered, the last application of the formula (for  $n = 6$ ), the minimum was produced by  $d_2 = 3$ . The second minimum (for  $n = 6 - 3$ ) was also produced for a coin of that denomination. Thus, the minimum-coin set for  $n = 6$  is two 3's. ■

**EXAMPLE 3** *Coin-collecting problem* Several coins are placed in cells of an  $n \times m$  board, no more than one coin per cell. A robot, located in the upper left cell of the board, needs to collect as many of the coins as possible and bring them to the bottom right cell. On each step, the robot can move either one cell to the right or one cell down from its current location. When the robot visits a cell with a coin, it always picks up that coin. Design an algorithm to find the maximum number of coins the robot can collect and a path it needs to follow to do this.

Let  $F(i, j)$  be the largest number of coins the robot can collect and bring to the cell  $(i, j)$  in the  $i$ th row and  $j$ th column of the board. It can reach this cell either from the adjacent cell  $(i - 1, j)$  above it or from the adjacent cell  $(i, j - 1)$  to the left of it. The largest numbers of coins that can be brought to these cells are  $F(i - 1, j)$  and  $F(i, j - 1)$ , respectively. Of course, there are no adjacent cells

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.

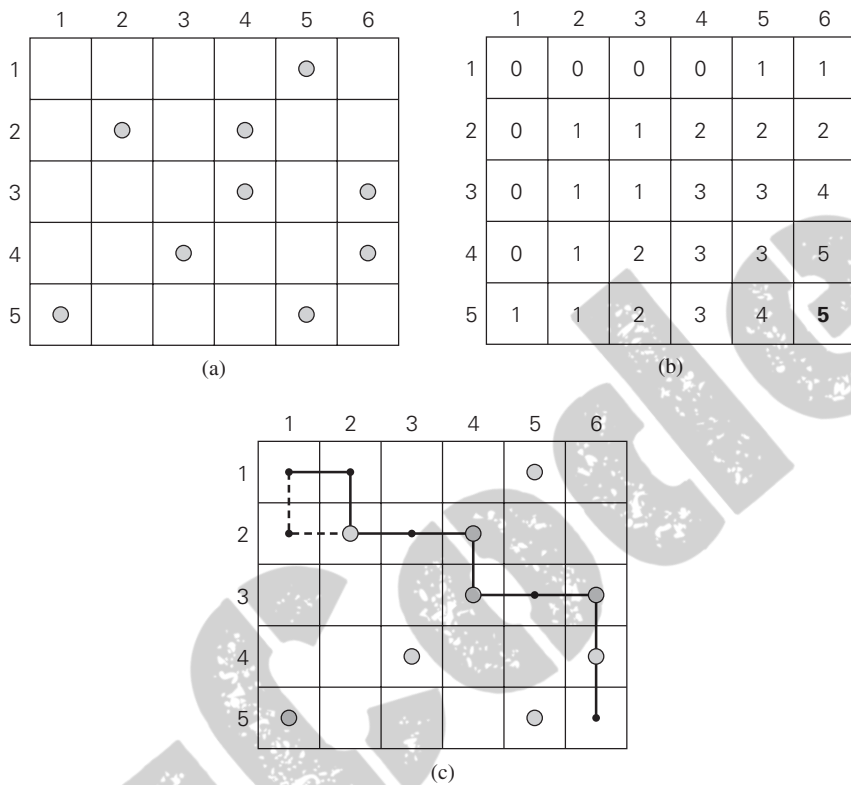
Using these formulas, we can fill in the  $n \times m$  table of  $F(i, j)$  values either row by row or column by column, as is typical for dynamic programming algorithms involving two-dimensional tables.

**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.

## 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
$w_1 = 2, v_1 = 12$ $w_2 = 1, v_2 = 10$ $w_3 = 3, v_3 = 20$ $w_4 = 2, v_4 = 15$		0	0	0	0	0	0	0
	1	0	0	12	12	12	12	12
	2	0	10	12	22	22	22	22
	3	0	10	12	22	30	32	32
	4	0	10	15	25	30	<b>37</b>	

**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	—	—	—	—	<b>37</b>

**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.

## 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

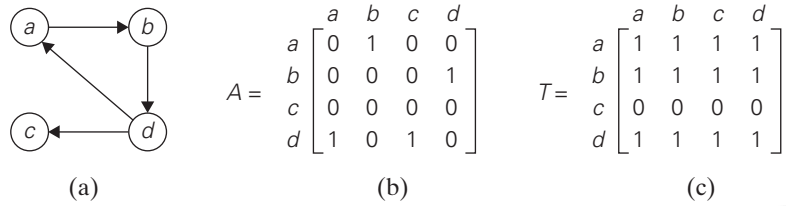
Recall that the adjacency matrix  $A = \{a_{ij}\}$  of a directed graph is the boolean matrix that has 1 in its  $i$ th row and  $j$ th column if and only if there is a directed edge from the  $i$ th vertex to the  $j$ th vertex. We may also be interested in a matrix containing the information about 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,  $r_{kj}^{(k-1)} = 1$ ).

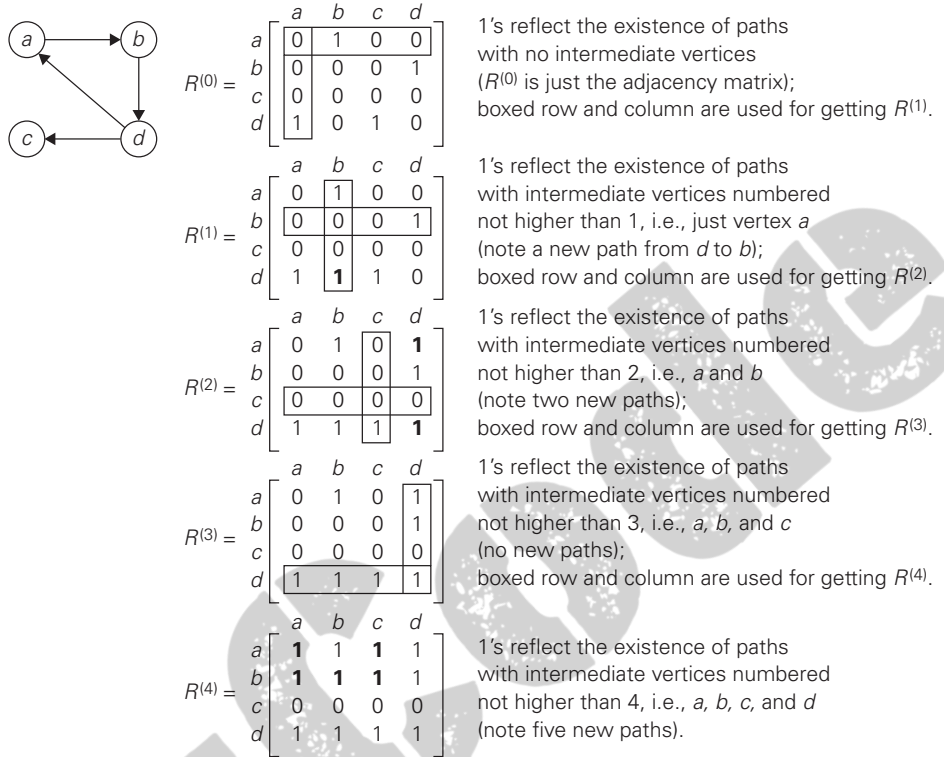
What we have just proved is that if  $r_{ij}^{(k)} = 1$ , then either  $r_{ij}^{(k-1)} = 1$  or both  $r_{ik}^{(k-1)} = 1$  and  $r_{kj}^{(k-1)} = 1$ . It is easy to see that the converse of this assertion is also true. Thus, we have the following formula for generating the elements of matrix  $R^{(k)}$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

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

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.<sup>1</sup> It is applicable to both undirected and directed weighted graphs provided



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  $k$ :

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

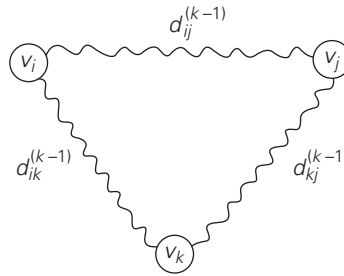
We can partition all such paths into two disjoint subsets: those that do not use the  $k$ th vertex  $v_k$  as intermediate and those that do. Since the paths of the first subset have their intermediate vertices numbered not higher than  $k - 1$ , the shortest of them is, by definition of our matrices,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

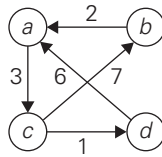
**for**  $i \leftarrow 1$  **to**  $n$  **do**

**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.



$$D^{(0)} = \begin{bmatrix} a & b & c & d \\ a & 0 & \infty & 3 & \infty \\ b & 2 & 0 & \infty & \infty \\ c & \infty & 7 & 0 & 1 \\ d & 6 & \infty & \infty & 0 \end{bmatrix}$$

Lengths of the shortest paths with no intermediate vertices ( $D^{(0)}$  is simply the weight matrix).

$$D^{(1)} = \begin{bmatrix} a & b & c & d \\ a & 0 & \infty & 3 & \infty \\ b & 2 & 0 & \mathbf{5} & \infty \\ c & \infty & 7 & 0 & 1 \\ d & 6 & \infty & \mathbf{9} & 0 \end{bmatrix}$$

Lengths of the shortest paths with intermediate vertices numbered not higher than 1, i.e., just  $a$  (note two new shortest paths from  $b$  to  $c$  and from  $d$  to  $c$ ).

$$D^{(2)} = \begin{bmatrix} a & b & c & d \\ a & 0 & \infty & 3 & \infty \\ b & 2 & 0 & 5 & \infty \\ c & \mathbf{9} & 7 & 0 & 1 \\ d & 6 & \infty & 9 & 0 \end{bmatrix}$$

Lengths of the shortest paths with intermediate vertices numbered not higher than 2, i.e.,  $a$  and  $b$  (note a new shortest path from  $c$  to  $a$ ).

$$D^{(3)} = \begin{bmatrix} a & b & c & d \\ a & 0 & \mathbf{10} & 3 & \mathbf{4} \\ b & 2 & 0 & 5 & \mathbf{6} \\ c & 9 & 7 & 0 & 1 \\ d & 6 & \mathbf{16} & 9 & 0 \end{bmatrix}$$

Lengths of the shortest paths with intermediate vertices numbered not higher than 3, i.e.,  $a$ ,  $b$ , and  $c$  (note four new shortest paths from  $a$  to  $b$ , from  $a$  to  $d$ , from  $b$  to  $d$ , and from  $d$  to  $b$ ).

$$D^{(4)} = \begin{bmatrix} a & b & c & d \\ a & 0 & 10 & 3 & 4 \\ b & 2 & 0 & 5 & 6 \\ c & \mathbf{7} & 7 & 0 & 1 \\ d & 6 & 16 & 9 & 0 \end{bmatrix}$$

Lengths of the shortest paths with intermediate vertices numbered not higher than 4, i.e.,  $a$ ,  $b$ ,  $c$ , and  $d$  (note a new shortest path from  $c$  to  $a$ ).

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

## 9.1 Prim's Algorithm

The following problem arises naturally in many practical situations: given  $n$  points, connect them in the cheapest possible way so that there will be a path between every pair of points. It has direct applications to the design of all kinds of networks—including communication, computer, transportation, and electrical—by providing the cheapest way to achieve connectivity. It identifies clusters of points in data sets. It has been used for classification purposes in archeology, biology, sociology, and other sciences. It is also helpful for constructing approximate solutions to more difficult problems such as the traveling salesman problem (see Section 12.3).

We can represent the points given by vertices of a graph, possible connections by the graph's edges, and the connection costs by the edge weights. Then the question can be posed as the minimum spanning tree problem, defined formally as follows.

**DEFINITION** A *spanning tree* of an undirected connected graph is its connected acyclic subgraph (i.e., a tree) that contains all the vertices of the graph. If such a graph has weights assigned to its edges, a *minimum spanning tree* is its spanning tree of the smallest weight, where the *weight* of a tree is defined as the sum of the weights on all its edges. The *minimum spanning tree problem* is the problem of finding a minimum spanning tree for a given weighted connected graph.

Figure 9.2 presents a simple example illustrating these notions.

If we were to try constructing a minimum spanning tree by exhaustive search, we would face two serious obstacles. First, the number of spanning trees grows exponentially with the graph size (at least for dense graphs). Second, generating all spanning trees for a given graph is not easy; in fact, it is more difficult than finding a *minimum* spanning tree for a weighted graph by using one of several efficient algorithms available for this problem. In this section, we outline *Prim's algorithm*, which goes back to at least 1957<sup>1</sup> [Pri57].

Prim's algorithm constructs a minimum spanning tree through a sequence of expanding subtrees. The initial subtree in such a sequence consists of a single vertex selected arbitrarily from the set  $V$  of the graph's vertices. On each iteration, the algorithm expands the current tree in the greedy manner by simply attaching to it the nearest vertex not in that tree. (By the nearest vertex, we mean a vertex not in the tree connected to a vertex in the tree by an edge of the smallest weight. Ties can be broken arbitrarily.) The algorithm stops after all the graph's vertices have been included in the tree being constructed. Since the algorithm expands a tree by exactly one vertex on each of its iterations, the total number of such iterations is  $n - 1$ , where  $n$  is the number of vertices in the graph. The tree generated by the algorithm is obtained as the set of edges used for the tree expansions.

Here is pseudocode of this algorithm.

**ALGORITHM** *Prim*( $G$ )

```
//Prim's algorithm for constructing a minimum spanning tree
//Input: A weighted connected graph  $G = \langle V, E \rangle$ 
//Output:  $E_T$ , the set of edges composing a minimum spanning tree of  $G$ 
 $V_T \leftarrow \{v_0\}$  //the set of tree vertices can be initialized with any vertex
 $E_T \leftarrow \emptyset$ 
for  $i \leftarrow 1$  to  $|V| - 1$  do
    find a minimum-weight edge  $e^* = (v^*, u^*)$  among all the edges  $(v, u)$ 
    such that  $v$  is in  $V_T$  and  $u$  is in  $V - V_T$ 
     $V_T \leftarrow V_T \cup \{u^*\}$ 
     $E_T \leftarrow E_T \cup \{e^*\}$ 
return  $E_T$ 
```

The nature of Prim's algorithm makes it necessary to provide each vertex not in the current tree with the information about the shortest edge connecting the vertex to a tree vertex. We can provide such information by attaching two labels to a vertex: the name of the nearest tree vertex and the length (the weight) of the corresponding edge. Vertices that are not adjacent to any of the tree vertices can be given the  $\infty$  label indicating their "infinite" distance to the tree vertices and a null label for the name of the nearest tree vertex. (Alternatively, we can split the vertices that are not in the tree into two sets, the "fringe" and the "unseen." The fringe contains only the vertices that are not in the tree but are adjacent to at least one tree vertex. These are the candidates from which the next tree vertex is selected. The unseen vertices are all the other vertices of the graph, called "unseen" because they are yet to be affected by the algorithm.) With such labels, finding the next vertex to be added to the current tree  $T = \langle V_T, E_T \rangle$  becomes a simple task of finding a vertex with the smallest distance label in the set  $V - V_T$ . Ties can be broken arbitrarily.

After we have identified a vertex  $u^*$  to be added to the tree, we need to perform two operations:

- Move  $u^*$  from the set  $V - V_T$  to the set of tree vertices  $V_T$ .
- For each remaining vertex  $u$  in  $V - V_T$  that is connected to  $u^*$  by a shorter edge than the  $u$ 's current distance label, update its labels by  $u^*$  and the weight of the edge between  $u^*$  and  $u$ , respectively.<sup>2</sup>

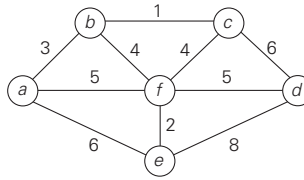
Figure 9.3 demonstrates the application of Prim's algorithm to a specific graph.

Does Prim's algorithm always yield a minimum spanning tree? The answer to this question is yes. Let us prove by induction that each of the subtrees  $T_i$ ,  $i = 0, \dots, n - 1$ , generated by Prim's algorithm is a part (i.e., a subgraph) of some minimum spanning tree. (This immediately implies, of course, that the last tree in the sequence,  $T_{n-1}$ , is a minimum spanning tree itself because it contains all  $n$  vertices of the graph.) The basis of the induction is trivial, since  $T_0$  consists of a single vertex and hence must be a part of any minimum spanning tree. For the inductive step, let us assume that  $T_{i-1}$  is part of some minimum spanning tree  $T$ . We need to prove that  $T_i$ , generated from  $T_{i-1}$  by Prim's algorithm, is also a part of a minimum spanning tree. We prove this by contradiction by assuming that no minimum spanning tree of the graph can contain  $T_i$ . Let  $e_i = (v, u)$  be the minimum weight edge from a vertex in  $T_{i-1}$  to a vertex not in  $T_{i-1}$  used by Prim's algorithm to expand  $T_{i-1}$  to  $T_i$ . By our assumption,  $e_i$  cannot belong to any minimum spanning tree, including  $T$ . Therefore, if we add  $e_i$  to  $T$ , a cycle must be formed (Figure 9.4).

In addition to edge  $e_i = (v, u)$ , this cycle must contain another edge  $(v', u')$  connecting a vertex  $v' \in T_{i-1}$  to a vertex  $u'$  that is not in  $T_{i-1}$ . (It is possible that  $v'$  coincides with  $v$  or  $u'$  coincides with  $u$  but not both.) If we now delete the edge  $(v', u')$  from this cycle, we will obtain another spanning tree of the entire graph whose weight is less than or equal to the weight of  $T$  since the weight of  $e_i$  is less than or equal to the weight of  $(v', u')$ . Hence, this spanning tree is a minimum spanning tree, which contradicts the assumption that no minimum spanning tree contains  $T_i$ . This completes the correctness proof of Prim's algorithm.

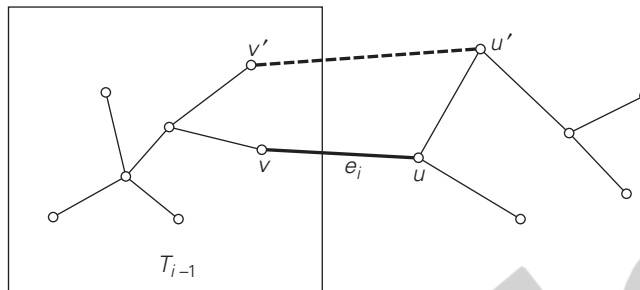
How efficient is Prim's algorithm? The answer depends on the data structures chosen for the graph itself and for the priority queue of the set  $V - V_T$  whose vertex priorities are the distances to the nearest tree vertices. (You may want to take another look at the example in Figure 9.3 to see that the set  $V - V_T$  indeed operates as a priority queue.) In particular, if a graph is represented by its weight matrix and the priority queue is implemented as an unordered array, the algorithm's running time will be in  $\Theta(|V|^2)$ . Indeed, on each of the  $|V| - 1$  iterations, the array implementing the priority queue is traversed to find and delete the minimum and then to update, if necessary, the priorities of the remaining vertices.

We can also implement the priority queue as a **min-heap**. A min-heap is a mirror image of the heap structure discussed in Section 6.4. (In fact, it can be implemented by constructing a heap after negating all the key values given.) Namely, a min-heap is a complete binary tree in which every element is less than or equal



Tree vertices	Remaining vertices	Illustration
$a(-, -)$	<b><math>b(a, 3)</math></b> $c(-, \infty)$ $d(-, \infty)$ $e(a, 6)$ $f(a, 5)$	
$b(a, 3)$	<b><math>c(b, 1)</math></b> $d(-, \infty)$ $e(a, 6)$ $f(b, 4)$	
$c(b, 1)$	$d(c, 6)$ $e(a, 6)$ <b><math>f(b, 4)</math></b>	
$f(b, 4)$	$d(f, 5)$ <b><math>e(f, 2)</math></b>	
$e(f, 2)$	<b><math>d(f, 5)</math></b>	
$d(f, 5)$		

**FIGURE 9.3** Application of Prim's algorithm. The parenthesized labels of a vertex in the middle column indicate the nearest tree vertex and edge weight; selected vertices and edges are shown in bold.



**FIGURE 9.4** Correctness proof of Prim's algorithm.

to its children. All the principal properties of heaps remain valid for min-heaps, with some obvious modifications. For example, the root of a min-heap contains the smallest rather than the largest element. Deletion of the smallest element from and insertion of a new element into a min-heap of size  $n$  are  $O(\log n)$  operations, and so is the operation of changing an element's priority (see Problem 15 in this section's exercises).

If a graph is represented by its adjacency lists and the priority queue is implemented as a min-heap, the running time of the algorithm is in  $O(|E| \log |V|)$ . This is because the algorithm performs  $|V| - 1$  deletions of the smallest element and makes  $|E|$  verifications and, possibly, changes of an element's priority in a min-heap of size not exceeding  $|V|$ . Each of these operations, as noted earlier, is a  $O(\log |V|)$  operation. Hence, the running time of this implementation of Prim's algorithm is in

$$(|V| - 1 + |E|)O(\log |V|) = O(|E| \log |V|)$$

because, in a connected graph,  $|V| - 1 \leq |E|$ .

In the next section, you will find another greedy algorithm for the minimum spanning tree problem, which is "greedy" in a manner different from that of Prim's algorithm.



## 9.2 Kruskal's Algorithm

In the previous section, we considered the greedy algorithm that “grows” a minimum spanning tree through a greedy inclusion of the nearest vertex to the vertices already in the tree. Remarkably, there is another greedy algorithm for the minimum spanning tree problem that also always yields an optimal solution. It is named **Kruskal's algorithm** after Joseph Kruskal, who discovered this algorithm when he was a second-year graduate student [Kru56]. Kruskal's algorithm looks at a minimum spanning tree of a weighted connected graph  $G = (V, E)$  as an acyclic subgraph with  $|V| - 1$  edges for which the sum of the edge weights is the smallest. (It is not difficult to prove that such a subgraph must be a tree.) Consequently, the algorithm constructs a minimum spanning tree as an expanding sequence of subgraphs that are always acyclic but are not necessarily connected on the intermediate stages of the algorithm.

The algorithm begins by sorting the graph's edges in nondecreasing order of their weights. Then, starting with the empty subgraph, it scans this sorted list, adding the next edge on the list to the current subgraph if such an inclusion does not create a cycle and simply skipping the edge otherwise.

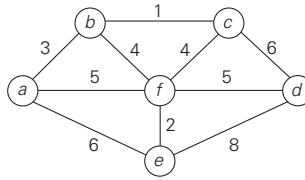
### ALGORITHM *Kruskal*( $G$ )

```
//Kruskal's algorithm for constructing a minimum spanning tree
//Input: A weighted connected graph  $G = (V, E)$ 
//Output:  $E_T$ , the set of edges composing a minimum spanning tree of  $G$ 
sort  $E$  in nondecreasing order of the edge weights  $w(e_{i_1}) \leq \dots \leq w(e_{i_{|E|}})$ 
 $E_T \leftarrow \emptyset$ ;  $ecounter \leftarrow 0$  //initialize the set of tree edges and its size
 $k \leftarrow 0$  //initialize the number of processed edges
while  $ecounter < |V| - 1$  do
     $k \leftarrow k + 1$ 
    if  $E_T \cup \{e_{i_k}\}$  is acyclic
         $E_T \leftarrow E_T \cup \{e_{i_k}\}$ ;  $ecounter \leftarrow ecounter + 1$ 
return  $E_T$ 
```

The correctness of Kruskal's algorithm can be proved by repeating the essential steps of the proof of Prim's algorithm given in the previous section. The fact that  $E_T$  is actually a tree in Prim's algorithm but generally just an acyclic subgraph in Kruskal's algorithm turns out to be an obstacle that can be overcome.

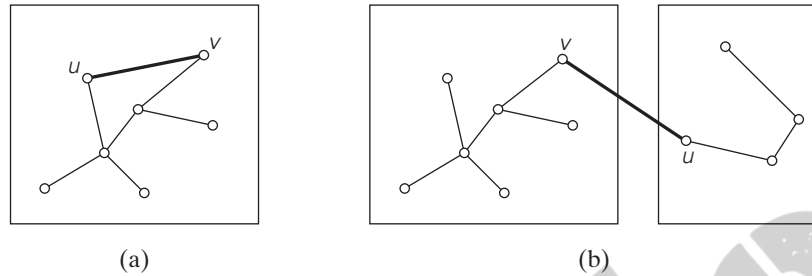
Figure 9.5 demonstrates the application of Kruskal's algorithm to the same graph we used for illustrating Prim's algorithm in Section 9.1. As you trace the algorithm's operations, note the disconnectedness of some of the intermediate subgraphs.

Applying Prim's and Kruskal's algorithms to the same small graph by hand may create the impression that the latter is simpler than the former. This impression is wrong because, on each of its iterations, Kruskal's algorithm has to check whether the addition of the next edge to the edges already selected would create a



Tree edges	Sorted list of edges										Illustration
	<b>bc</b> 1	ef 2	ab 3	bf 4	cf 4	af 5	df 5	ae 6	cd 6	de 8	
bc 1	bc 1	<b>ef</b> 2	ab 3	bf 4	cf 4	af 5	df 5	ae 6	cd 6	de 8	
ef 2	bc 1	ef 2	<b>ab</b> 3	bf 4	cf 4	af 5	df 5	ae 6	cd 6	de 8	
ab 3	bc 1	ef 2	ab 3	<b>bf</b> 4	cf 4	af 5	df 5	ae 6	cd 6	de 8	
bf 4	bc 1	ef 2	ab 3	bf 4	cf 4	af 5	<b>df</b> 5	ae 6	cd 6	de 8	
df 5											

**FIGURE 9.5** Application of Kruskal's algorithm. Selected edges are shown in bold.



**FIGURE 9.6** New edge connecting two vertices may (a) or may not (b) create a cycle.

cycle. It is not difficult to see that a new cycle is created if and only if the new edge connects two vertices already connected by a path, i.e., if and only if the two vertices belong to the same connected component (Figure 9.6). Note also that each connected component of a subgraph generated by Kruskal's algorithm is a tree because it has no cycles.

In view of these observations, it is convenient to use a slightly different interpretation of Kruskal's algorithm. We can consider the algorithm's operations as a progression through a series of forests containing *all* the vertices of a given graph and *some* of its edges. The initial forest consists of  $|V|$  trivial trees, each comprising a single vertex of the graph. The final forest consists of a single tree, which is a minimum spanning tree of the graph. On each iteration, the algorithm takes the next edge  $(u, v)$  from the sorted list of the graph's edges, finds the trees containing the vertices  $u$  and  $v$ , and, if these trees are not the same, unites them in a larger tree by adding the edge  $(u, v)$ .

Fortunately, there are efficient algorithms for doing so, including the crucial check for whether two vertices belong to the same tree. They are called **union-find** algorithms. We discuss them in the following subsection. With an efficient union-find algorithm, the running time of Kruskal's algorithm will be dominated by the time needed for sorting the edge weights of a given graph. Hence, with an efficient sorting algorithm, the time efficiency of Kruskal's algorithm will be in  $O(|E| \log |E|)$ .

## Disjoint Subsets and Union-Find Algorithms

Kruskal's algorithm is one of a number of applications that require a dynamic partition of some  $n$  element set  $S$  into a collection of disjoint subsets  $S_1, S_2, \dots, S_k$ . After being initialized as a collection of  $n$  one-element subsets, each containing a different element of  $S$ , the collection is subjected to a sequence of intermixed union and find operations. (Note that the number of union operations in any such sequence must be bounded above by  $n - 1$  because each union increases a subset's size at least by 1 and there are only  $n$  elements in the entire set  $S$ .) Thus, we are

dealing here with an abstract data type of a collection of disjoint subsets of a finite set with the following operations:

*makeset*( $x$ ) creates a one-element set  $\{x\}$ . It is assumed that this operation can be applied to each of the elements of set  $S$  only once.

*find*( $x$ ) returns a subset containing  $x$ .

*union*( $x, y$ ) constructs the union of the disjoint subsets  $S_x$  and  $S_y$  containing  $x$  and  $y$ , respectively, and adds it to the collection to replace  $S_x$  and  $S_y$ , which are deleted from it.

For example, let  $S = \{1, 2, 3, 4, 5, 6\}$ . Then *makeset*( $i$ ) creates the set  $\{i\}$  and applying this operation six times initializes the structure to the collection of six singleton sets:

$\{1\}, \{2\}, \{3\}, \{4\}, \{5\}, \{6\}.$

Performing *union*(1, 4) and *union*(5, 2) yields

$\{1, 4\}, \{5, 2\}, \{3\}, \{6\},$

and, if followed by *union*(4, 5) and then by *union*(3, 6), we end up with the disjoint subsets

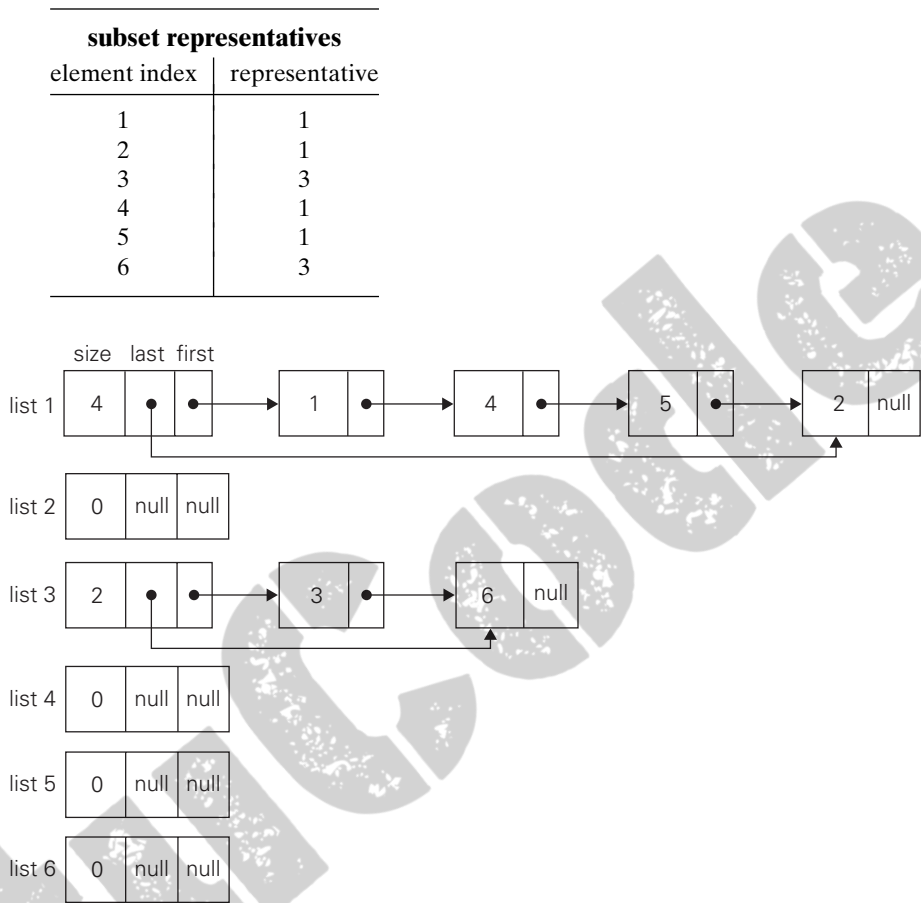
$\{1, 4, 5, 2\}, \{3, 6\}.$

Most implementations of this abstract data type use one element from each of the disjoint subsets in a collection as that subset's **representative**. Some implementations do not impose any specific constraints on such a representative; others do so by requiring, say, the smallest element of each subset to be used as the subset's representative. Also, it is usually assumed that set elements are (or can be mapped into) integers.

There are two principal alternatives for implementing this data structure. The first one, called the **quick find**, optimizes the time efficiency of the find operation; the second one, called the **quick union**, optimizes the union operation.

The quick find uses an array indexed by the elements of the underlying set  $S$ ; the array's values indicate the representatives of the subsets containing those elements. Each subset is implemented as a linked list whose header contains the pointers to the first and last elements of the list along with the number of elements in the list (see Figure 9.7 for an example).

Under this scheme, the implementation of *makeset*( $x$ ) requires assigning the corresponding element in the representative array to  $x$  and initializing the corresponding linked list to a single node with the  $x$  value. The time efficiency of this operation is obviously in  $\Theta(1)$ , and hence the initialization of  $n$  singleton subsets is in  $\Theta(n)$ . The efficiency of *find*( $x$ ) is also in  $\Theta(1)$ : all we need to do is to retrieve the  $x$ 's representative in the representative array. Executing *union*( $x, y$ ) takes longer. A straightforward solution would simply append the  $y$ 's list to the end of the  $x$ 's list, update the information about their representative for all the elements in the



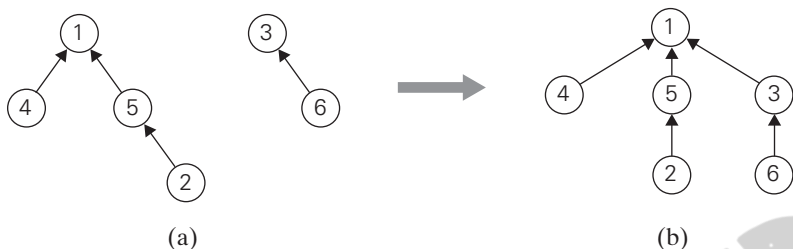
**FIGURE 9.7** Linked-list representation of subsets {1, 4, 5, 2} and {3, 6} obtained by quick find after performing  $\text{union}(1, 4)$ ,  $\text{union}(5, 2)$ ,  $\text{union}(4, 5)$ , and  $\text{union}(3, 6)$ . The lists of size 0 are considered deleted from the collection.

y list, and then delete the y's list from the collection. It is easy to verify, however, that with this algorithm the sequence of union operations

$$\text{union}(2, 1), \text{union}(3, 2), \dots, \text{union}(i + 1, i), \dots, \text{union}(n, n - 1)$$

runs in  $\Theta(n^2)$  time, which is slow compared with several known alternatives.

A simple way to improve the overall efficiency of a sequence of union operations is to always append the shorter of the two lists to the longer one, with ties broken arbitrarily. Of course, the size of each list is assumed to be available by, say, storing the number of elements in the list's header. This modification is called the



**FIGURE 9.8** (a) Forest representation of subsets  $\{1, 4, 5, 2\}$  and  $\{3, 6\}$  used by quick union. (b) Result of  $\text{union}(5, 6)$ .

**union by size.** Though it does not improve the worst-case efficiency of a single application of the union operation (it is still in  $\Theta(n)$ ), the worst-case running time of any legitimate sequence of union-by-size operations turns out to be in  $O(n \log n)$ .<sup>3</sup>

Here is a proof of this assertion. Let  $a_i$  be an element of set  $S$  whose disjoint subsets we manipulate, and let  $A_i$  be the number of times  $a_i$ 's representative is updated in a sequence of union-by-size operations. How large can  $A_i$  get if set  $S$  has  $n$  elements? Each time  $a_i$ 's representative is updated,  $a_i$  must be in a smaller subset involved in computing the union whose size will be at least twice as large as the size of the subset containing  $a_i$ . Hence, when  $a_i$ 's representative is updated for the first time, the resulting set will have at least two elements; when it is updated for the second time, the resulting set will have at least four elements; and, in general, if it is updated  $A_i$  times, the resulting set will have at least  $2^{A_i}$  elements. Since the entire set  $S$  has  $n$  elements,  $2^{A_i} \leq n$  and hence  $A_i \leq \log_2 n$ . Therefore, the total number of possible updates of the representatives for all  $n$  elements in  $S$  will not exceed  $n \log_2 n$ .

Thus, for union by size, the time efficiency of a sequence of at most  $n - 1$  unions and  $m$  finds is in  $O(n \log n + m)$ .

The **quick union**—the second principal alternative for implementing disjoint subsets—represents each subset by a rooted tree. The nodes of the tree contain the subset's elements (one per node), with the root's element considered the subset's representative; the tree's edges are directed from children to their parents (Figure 9.8). In addition, a mapping of the set elements to their tree nodes—implemented, say, as an array of pointers—is maintained. This mapping is not shown in Figure 9.8 for the sake of simplicity.

For this implementation,  $\text{makeset}(x)$  requires the creation of a single-node tree, which is a  $\Theta(1)$  operation; hence, the initialization of  $n$  singleton subsets is in  $\Theta(n)$ . A  $\text{union}(x, y)$  is implemented by attaching the root of the  $y$ 's tree to the root of the  $x$ 's tree (and deleting the  $y$ 's tree from the collection by making the pointer to its root null). The time efficiency of this operation is clearly  $\Theta(1)$ . A  $\text{find}(x)$  is

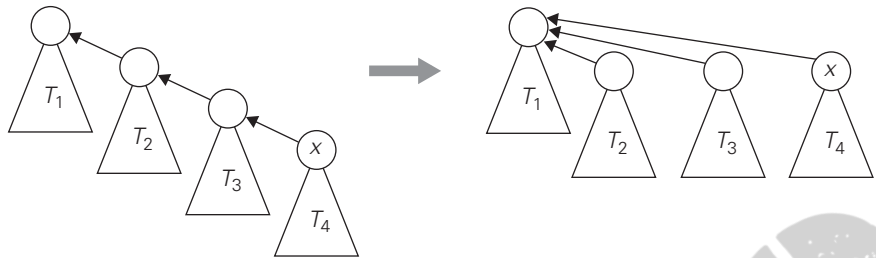


FIGURE 9.9 Path compression.

performed by following the pointer chain from the node containing  $x$  to the tree's root whose element is returned as the subset's representative. Accordingly, the time efficiency of a single find operation is in  $O(n)$  because a tree representing a subset can degenerate into a linked list with  $n$  nodes.

This time bound can be improved. The straightforward way for doing so is to always perform a union operation by attaching a smaller tree to the root of a larger one, with ties broken arbitrarily. The size of a tree can be measured either by the number of nodes (this version is called **union by size**) or by its height (this version is called **union by rank**). Of course, these options require storing, for each node of the tree, either the number of node descendants or the height of the subtree rooted at that node, respectively. One can easily prove that in either case the height of the tree will be logarithmic, making it possible to execute each find in  $O(\log n)$  time. Thus, for quick union, the time efficiency of a sequence of at most  $n - 1$  unions and  $m$  finds is in  $O(n + m \log n)$ .

In fact, an even better efficiency can be obtained by combining either variety of quick union with **path compression**. This modification makes every node encountered during the execution of a find operation point to the tree's root (Figure 9.9). According to a quite sophisticated analysis that goes beyond the level of this book (see [Tar84]), this and similar techniques improve the efficiency of a sequence of at most  $n - 1$  unions and  $m$  finds to only slightly worse than linear.

## 9.3 Dijkstra's Algorithm

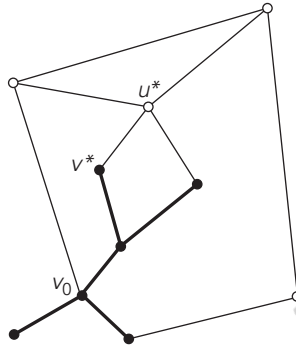
In this section, we consider the *single-source shortest-paths problem*: for a given vertex called the *source* in a weighted connected graph, find shortest paths to all its other vertices. It is important to stress that we are not interested here in a single shortest path that starts at the source and visits all the other vertices. This would have been a much more difficult problem (actually, a version of the traveling salesman problem introduced in Section 3.4 and discussed again later in the book). The single-source shortest-paths problem asks for a family of paths, each leading from the source to a different vertex in the graph, though some paths may, of course, have edges in common.

A variety of practical applications of the shortest-paths problem have made the problem a very popular object of study. The obvious but probably most widely used applications are transportation planning and packet routing in communication networks, including the Internet. Multitudes of less obvious applications include finding shortest paths in social networks, speech recognition, document formatting, robotics, compilers, and airline crew scheduling. In the world of entertainment, one can mention pathfinding in video games and finding best solutions to puzzles using their state-space graphs (see Section 6.6 for a very simple example of the latter).

There are several well-known algorithms for finding shortest paths, including Floyd's algorithm for the more general all-pairs shortest-paths problem discussed in Chapter 8. Here, we consider the best-known algorithm for the single-source shortest-paths problem, called *Dijkstra's algorithm*.<sup>4</sup> This algorithm is applicable to undirected and directed graphs with nonnegative weights only. Since in most applications this condition is satisfied, the limitation has not impaired the popularity of Dijkstra's algorithm.

Dijkstra's algorithm finds the shortest paths to a graph's vertices in order of their distance from a given source. First, it finds the shortest path from the source





**FIGURE 9.10** Idea of Dijkstra's algorithm. The subtree of the shortest paths already found is shown in bold. The next nearest to the source  $v_0$  vertex,  $u^*$ , is selected by comparing the lengths of the subtree's paths increased by the distances to vertices adjacent to the subtree's vertices.

to a vertex nearest to it, then to a second nearest, and so on. In general, before its  $i$ th iteration commences, the algorithm has already identified the shortest paths to  $i - 1$  other vertices nearest to the source. These vertices, the source, and the edges of the shortest paths leading to them from the source form a subtree  $T_i$  of the given graph (Figure 9.10). Since all the edge weights are nonnegative, the next vertex nearest to the source can be found among the vertices adjacent to the vertices of  $T_i$ . The set of vertices adjacent to the vertices in  $T_i$  can be referred to as “fringe vertices”; they are the candidates from which Dijkstra's algorithm selects the next vertex nearest to the source. (Actually, all the other vertices can be treated as fringe vertices connected to tree vertices by edges of infinitely large weights.) To identify the  $i$ th nearest vertex, the algorithm computes, for every fringe vertex  $u$ , the sum of the distance to the nearest tree vertex  $v$  (given by the weight of the edge  $(v, u)$ ) and the length  $d_v$  of the shortest path from the source to  $v$  (previously determined by the algorithm) and then selects the vertex with the smallest such sum. The fact that it suffices to compare the lengths of such special paths is the central insight of Dijkstra's algorithm.

To facilitate the algorithm's operations, we label each vertex with two labels. The numeric label  $d$  indicates the length of the shortest path from the source to this vertex found by the algorithm so far; when a vertex is added to the tree,  $d$  indicates the length of the shortest path from the source to that vertex. The other label indicates the name of the next-to-last vertex on such a path, i.e., the parent of the vertex in the tree being constructed. (It can be left unspecified for the source  $s$  and vertices that are adjacent to none of the current tree vertices.) With such labeling, finding the next nearest vertex  $u^*$  becomes a simple task of finding a fringe vertex with the smallest  $d$  value. Ties can be broken arbitrarily.

After we have identified a vertex  $u^*$  to be added to the tree, we need to perform two operations:

- Move  $u^*$  from the fringe to the set of tree vertices.
- For each remaining fringe vertex  $u$  that is connected to  $u^*$  by an edge of weight  $w(u^*, u)$  such that  $d_{u^*} + w(u^*, u) < d_u$ , update the labels of  $u$  by  $u^*$  and  $d_{u^*} + w(u^*, u)$ , respectively.

Figure 9.11 demonstrates the application of Dijkstra's algorithm to a specific graph.

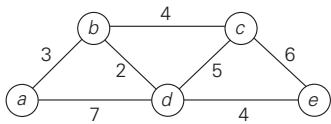
The labeling and mechanics of Dijkstra's algorithm are quite similar to those used by Prim's algorithm (see Section 9.1). Both of them construct an expanding subtree of vertices by selecting the next vertex from the priority queue of the remaining vertices. It is important not to mix them up, however. They solve different problems and therefore operate with priorities computed in a different manner: Dijkstra's algorithm compares path lengths and therefore must add edge weights, while Prim's algorithm compares the edge weights as given.

Now we can give pseudocode of Dijkstra's algorithm. It is spelled out—in more detail than Prim's algorithm was in Section 9.1—in terms of explicit operations on two sets of labeled vertices: the set  $V_T$  of vertices for which a shortest path has already been found and the priority queue  $Q$  of the fringe vertices. (Note that in the following pseudocode,  $V_T$  contains a given source vertex and the fringe contains the vertices adjacent to it *after* iteration 0 is completed.)

**ALGORITHM** *Dijkstra*( $G, s$ )

```
//Dijkstra's algorithm for single-source shortest paths
//Input: A weighted connected graph  $G = \langle V, E \rangle$  with nonnegative weights
//      and its vertex  $s$ 
//Output: The length  $d_v$  of a shortest path from  $s$  to  $v$ 
//      and its penultimate vertex  $p_v$  for every vertex  $v$  in  $V$ 
Initialize( $Q$ ) //initialize priority queue to empty
for every vertex  $v$  in  $V$ 
     $d_v \leftarrow \infty$ ;  $p_v \leftarrow \text{null}$ 
    Insert( $Q, v, d_v$ ) //initialize vertex priority in the priority queue
 $d_s \leftarrow 0$ ; Decrease( $Q, s, d_s$ ) //update priority of  $s$  with  $d_s$ 
 $V_T \leftarrow \emptyset$ 
for  $i \leftarrow 0$  to  $|V| - 1$  do
     $u^* \leftarrow \text{DeleteMin}(Q)$  //delete the minimum priority element
     $V_T \leftarrow V_T \cup \{u^*\}$ 
    for every vertex  $u$  in  $V - V_T$  that is adjacent to  $u^*$  do
        if  $d_{u^*} + w(u^*, u) < d_u$ 
             $d_u \leftarrow d_{u^*} + w(u^*, u)$ ;  $p_u \leftarrow u^*$ 
        Decrease( $Q, u, d_u$ )
```

The time efficiency of Dijkstra's algorithm depends on the data structures used for implementing the priority queue and for representing an input graph itself. For the reasons explained in the analysis of Prim's algorithm in Section 9.1, it is



Tree vertices	Remaining vertices	Illustration
$a(-, 0)$	<b><math>b(a, 3)</math></b> $c(-, \infty)$ $d(a, 7)$ $e(-, \infty)$	
$b(a, 3)$	$c(b, 3 + 4)$ <b><math>d(b, 3 + 2)</math></b> $e(-, \infty)$	
$d(b, 5)$	<b><math>c(b, 7)</math></b> $e(d, 5 + 4)$	
$c(b, 7)$	<b><math>e(d, 9)</math></b>	
$e(d, 9)$		

The shortest paths (identified by following nonnumeric labels backward from a destination vertex in the left column to the source) and their lengths (given by numeric labels of the tree vertices) are as follows:

- from  $a$  to  $b$  :  $a - b$  of length 3
- from  $a$  to  $d$  :  $a - b - d$  of length 5
- from  $a$  to  $c$  :  $a - b - c$  of length 7
- from  $a$  to  $e$  :  $a - b - d - e$  of length 9

**FIGURE 9.11** Application of Dijkstra’s algorithm. The next closest vertex is shown in bold.

in  $\Theta(|V|^2)$  for graphs represented by their weight matrix and the priority queue implemented as an unordered array. For graphs represented by their adjacency lists and the priority queue implemented as a min-heap, it is in  $O(|E| \log |V|)$ . A still better upper bound can be achieved for both Prim's and Dijkstra's algorithms if the priority queue is implemented using a sophisticated data structure called the ***Fibonacci heap*** (e.g., [Cor09]). However, its complexity and a considerable overhead make such an improvement primarily of theoretical value.

WUJICODLE

## 9.4 Huffman Trees and Codes

Suppose we have to encode a text that comprises symbols from some  $n$ -symbol alphabet by assigning to each of the text's symbols some sequence of bits called the **codeword**. For example, we can use a **fixed-length encoding** that assigns to each symbol a bit string of the same length  $m$  ( $m \geq \log_2 n$ ). This is exactly what the standard ASCII code does. One way of getting a coding scheme that yields a shorter bit string on the average is based on the old idea of assigning shorter codewords to more frequent symbols and longer codewords to less frequent symbols. This idea was used, in particular, in the telegraph code invented in the mid-19th century by Samuel Morse. In that code, frequent letters such as  $e$  ( $\cdot$ ) and  $a$  ( $\cdot -$ ) are assigned short sequences of dots and dashes while infrequent letters such as  $q$  ( $- - \cdot -$ ) and  $z$  ( $- - - \cdot$ ) have longer ones.

**Variable-length encoding**, which assigns codewords of different lengths to different symbols, introduces a problem that fixed-length encoding does not have. Namely, how can we tell how many bits of an encoded text represent the first (or, more generally, the  $i$ th) symbol? To avoid this complication, we can limit ourselves to the so-called **prefix-free** (or simply **prefix**) **codes**. In a prefix code, no codeword is a prefix of a codeword of another symbol. Hence, with such an encoding, we can simply scan a bit string until we get the first group of bits that is a codeword for some symbol, replace these bits by this symbol, and repeat this operation until the bit string's end is reached.

If we want to create a binary prefix code for some alphabet, it is natural to associate the alphabet's symbols with leaves of a binary tree in which all the left edges are labeled by 0 and all the right edges are labeled by 1. The codeword of a symbol can then be obtained by recording the labels on the simple path from the root to the symbol's leaf. Since there is no simple path to a leaf that continues to another leaf, no codeword can be a prefix of another codeword; hence, any such tree yields a prefix code.

Among the many trees that can be constructed in this manner for a given alphabet with known frequencies of the symbol occurrences, how can we construct a tree that would assign shorter bit strings to high-frequency symbols and longer ones to low-frequency symbols? It can be done by the following greedy algorithm, invented by David Huffman while he was a graduate student at MIT [Huf52].

### Huffman's algorithm

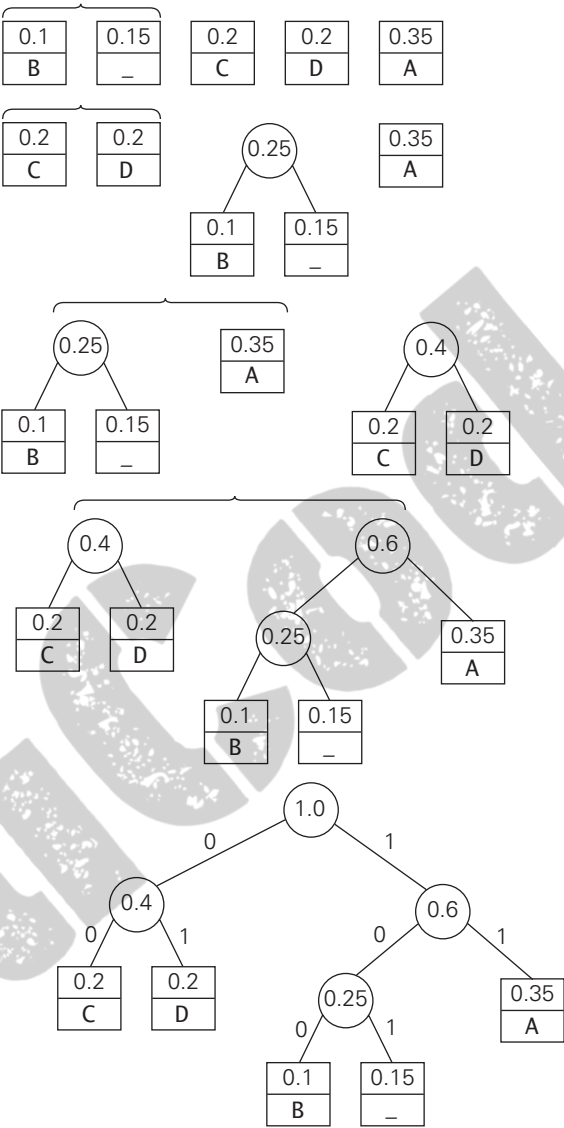
- Step 1** Initialize  $n$  one-node trees and label them with the symbols of the alphabet given. Record the frequency of each symbol in its tree's root to indicate the tree's **weight**. (More generally, the weight of a tree will be equal to the sum of the frequencies in the tree's leaves.)
- Step 2** Repeat the following operation until a single tree is obtained. Find two trees with the smallest weight (ties can be broken arbitrarily, but see Problem 2 in this section's exercises). Make them the left and right subtree of a new tree and record the sum of their weights in the root of the new tree as its weight.

A tree constructed by the above algorithm is called a **Huffman tree**. It defines—in the manner described above—a **Huffman code**.

**EXAMPLE** Consider the five-symbol alphabet {A, B, C, D, \_} with the following occurrence frequencies in a text made up of these symbols:

symbol	A	B	C	D	_
frequency	0.35	0.1	0.2	0.2	0.15

The Huffman tree construction for this input is shown in Figure 9.12.



**FIGURE 9.12** Example of constructing a Huffman coding tree.

The resulting codewords are as follows:

symbol	A	B	C	D	_
frequency	0.35	0.1	0.2	0.2	0.15
codeword	11	100	00	01	101

Hence, DAD is encoded as 011101, and 10011011011101 is decoded as BAD\_AD.

With the occurrence frequencies given and the codeword lengths obtained, the average number of bits per symbol in this code is

$$2 \cdot 0.35 + 3 \cdot 0.1 + 2 \cdot 0.2 + 2 \cdot 0.2 + 3 \cdot 0.15 = 2.25.$$

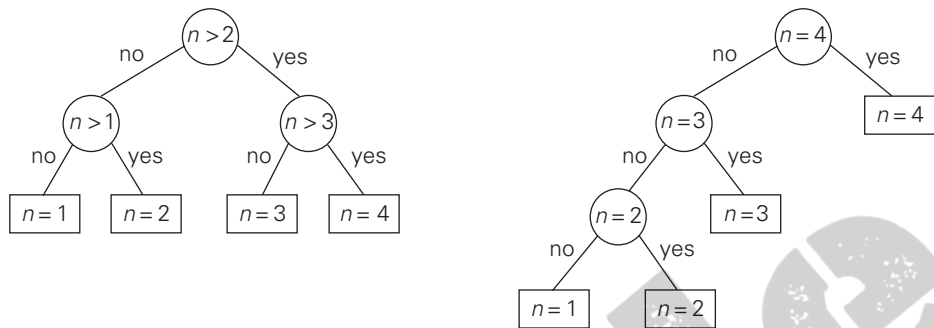
Had we used a fixed-length encoding for the same alphabet, we would have to use at least 3 bits per each symbol. Thus, for this toy example, Huffman's code achieves the **compression ratio**—a standard measure of a compression algorithm's effectiveness—of  $(3 - 2.25)/3 \cdot 100\% = 25\%$ . In other words, Huffman's encoding of the text will use 25% less memory than its fixed-length encoding. (Extensive experiments with Huffman codes have shown that the compression ratio for this scheme typically falls between 20% and 80%, depending on the characteristics of the text being compressed.) ■

Huffman's encoding is one of the most important file-compression methods. In addition to its simplicity and versatility, it yields an optimal, i.e., minimal-length, encoding (provided the frequencies of symbol occurrences are independent and known in advance). The simplest version of Huffman compression calls, in fact, for a preliminary scanning of a given text to count the frequencies of symbol occurrences in it. Then these frequencies are used to construct a Huffman coding tree and encode the text as described above. This scheme makes it necessary, however, to include the coding table into the encoded text to make its decoding possible. This drawback can be overcome by using **dynamic Huffman encoding**, in which the coding tree is updated each time a new symbol is read from the source text. Further, modern alternatives such as **Lempel-Ziv** algorithms (e.g., [Say05]) assign codewords not to individual symbols but to strings of symbols, allowing them to achieve better and more robust compressions in many applications.

It is important to note that applications of Huffman's algorithm are not limited to data compression. Suppose we have  $n$  positive numbers  $w_1, w_2, \dots, w_n$  that have to be assigned to  $n$  leaves of a binary tree, one per node. If we define the **weighted path length** as the sum  $\sum_{i=1}^n l_i w_i$ , where  $l_i$  is the length of the simple path from the root to the  $i$ th leaf, how can we construct a binary tree with minimum weighted path length? It is this more general problem that Huffman's algorithm actually solves. (For the coding application,  $l_i$  and  $w_i$  are the length of the codeword and the frequency of the  $i$ th symbol, respectively.)

This problem arises in many situations involving decision making. Consider, for example, the game of guessing a chosen object from  $n$  possibilities (say, an integer between 1 and  $n$ ) by asking questions answerable by yes or no. Different strategies for playing this game can be modeled by **decision trees**<sup>5</sup> such as those depicted in Figure 9.13 for  $n = 4$ . The length of the simple path from the root to a leaf in such a tree is equal to the number of questions needed to get to the chosen number represented by the leaf. If number  $i$  is chosen with probability  $p_i$ , the sum





**FIGURE 9.13** Two decision trees for guessing an integer between 1 and 4.

$\sum_{i=1}^n l_i p_i$ , where  $l_i$  is the length of the path from the root to the  $i$ th leaf, indicates the average number of questions needed to “guess” the chosen number with a game strategy represented by its decision tree. If each of the numbers is chosen with the same probability of  $1/n$ , the best strategy is to successively eliminate half (or almost half) the candidates as binary search does. This may not be the case for arbitrary  $p_i$ ’s, however. For example, if  $n = 4$  and  $p_1 = 0.1$ ,  $p_2 = 0.2$ ,  $p_3 = 0.3$ , and  $p_4 = 0.4$ , the minimum weighted path tree is the rightmost one in Figure 9.13. Thus, we need Huffman’s algorithm to solve this problem in its general case.

Note that this is the second time we are encountering the problem of constructing an optimal binary tree. In Section 8.3, we discussed the problem of constructing an optimal binary search tree with positive numbers (the search probabilities) assigned to every node of the tree. In this section, given numbers are assigned just to leaves. The latter problem turns out to be easier: it can be solved by the greedy algorithm, whereas the former is solved by the more complicated dynamic programming algorithm.