

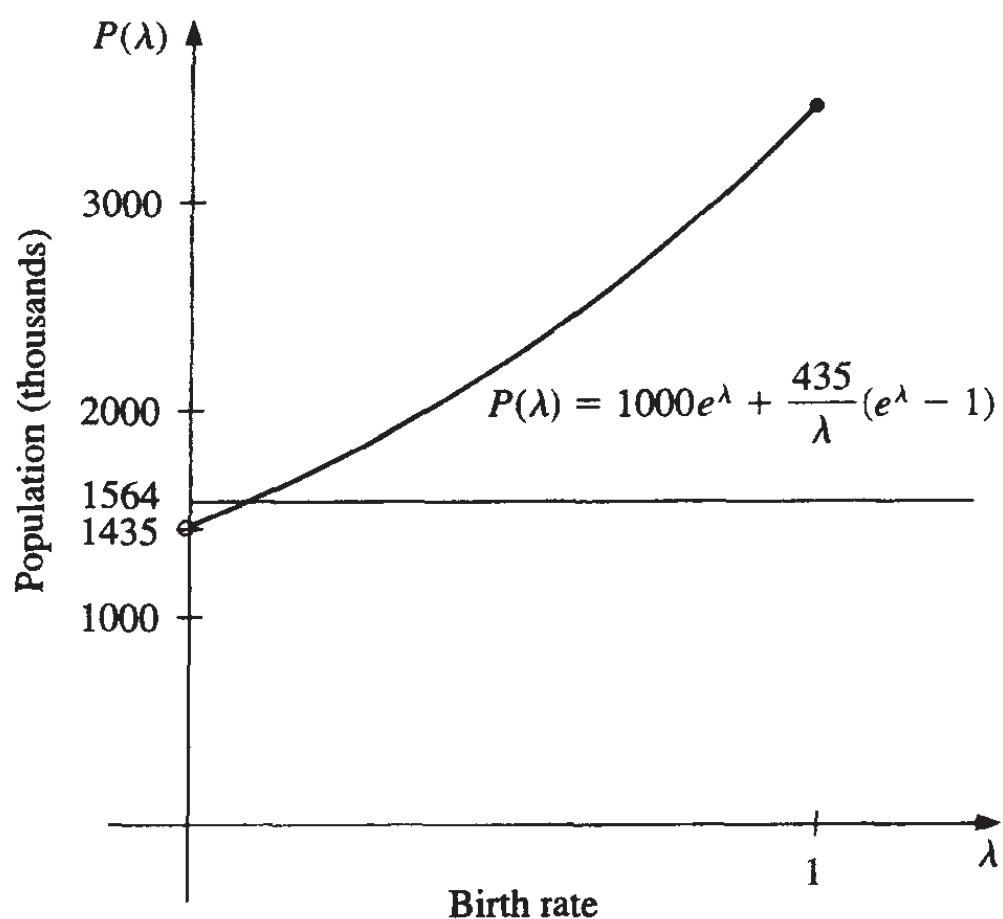
# Solutions of Equations in One Variable

■ ■ ■

The growth of a population can be modeled over short periods of time by assuming that the population grows continuously with time at a rate proportional to the number present at that time. If we let  $N(t)$  denote the number at time  $t$  and  $\lambda$  denote the constant birth rate of the population, then the population satisfies the differential equation

$$\frac{dN(t)}{dt} = \lambda N(t).$$

The solution to this equation is  $N(t) = N_0 e^{\lambda t}$ , where  $N_0$  denotes the initial population.



This exponential model is valid only when the population is isolated, with no immigration. If immigration is permitted at a constant rate  $\nu$ , then the differential equation becomes

$$\frac{dN(t)}{dt} = \lambda N(t) + \nu,$$

whose solution is

$$N(t) = N_0 e^{\lambda t} + \frac{\nu}{\lambda} (e^{\lambda t} - 1).$$

Suppose a certain population contains 1,000,000 individuals initially, that 435,000 individuals immigrate into the community in the first year, and that 1,564,000 individuals are present at the end of one year. To determine the birth rate of this population, we must solve for  $\lambda$  in the equation

$$1,564,000 = 1,000,000 e^{\lambda} + \frac{435,000}{\lambda} (e^{\lambda} - 1).$$

The numerical methods discussed in this chapter are used to approximate solutions of equations of this type, when the exact solutions cannot be obtained by algebraic methods. The solution to this particular problem is considered in Exercise 20 of Section 2.3.

## 2.1 The Bisection Method

In this chapter, we consider one of the most basic problems of numerical approximation, the *root-finding problem*. This process involves finding a *root*, or solution, of an equation of the form  $f(x) = 0$ , for a given function  $f$ . A root of this equation is also called a *zero* of the function  $f$ . The problem of finding an approximation to the root of an equation can be traced back at least as far as 1700 B.C. A cuneiform table in the Yale Babylonian Collection dating from that period gives a sexagesimal (base-60) number equivalent to 1.414222 as an approximation to  $\sqrt{2}$ , a result that is accurate to within  $10^{-5}$ . This approximation can be found by applying a technique described in Exercise 19 of Section 2.2.

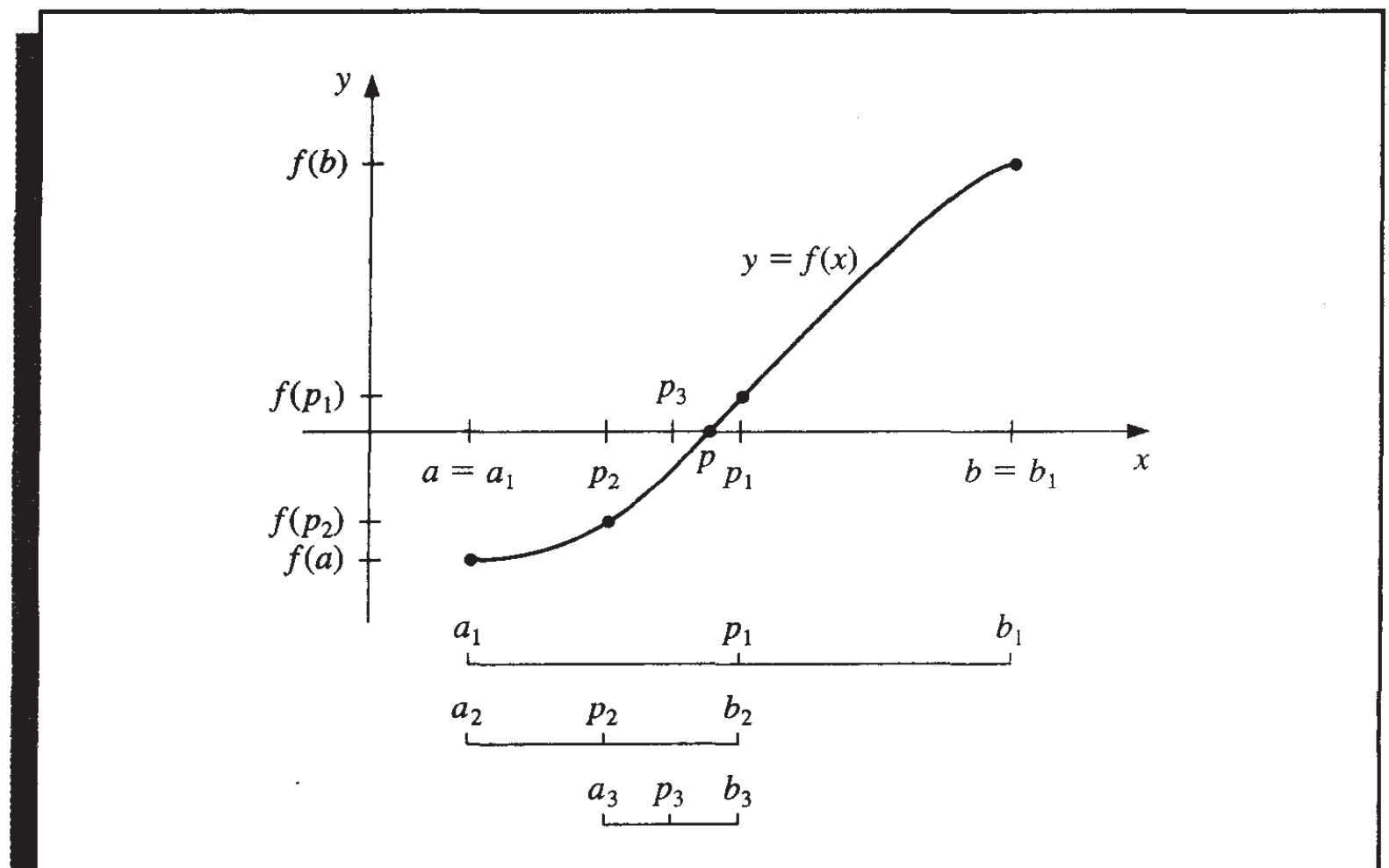
The first technique, based on the Intermediate Value Theorem, is called the **Bisection**, or **Binary-search, method**. Suppose  $f$  is a continuous function defined on the interval  $[a, b]$ , with  $f(a)$  and  $f(b)$  of opposite sign. By the Intermediate Value Theorem, there exists a number  $p$  in  $(a, b)$  with  $f(p) = 0$ . Although the procedure will work when there is more than one root in the interval  $(a, b)$ , we assume for simplicity that the root in this interval is unique. The method calls for a repeated halving of subintervals of  $[a, b]$  and, at each step, locating the half containing  $p$ .

To begin, set  $a_1 = a$  and  $b_1 = b$ , and let  $p_1$  be the midpoint of  $[a, b]$ ; that is,

$$p_1 = a_1 + \frac{b_1 - a_1}{2} = \frac{a_1 + b_1}{2}.$$

If  $f(p_1) = 0$ , then  $p = p_1$ , and we are done. If  $f(p_1) \neq 0$ , then  $f(p_1)$  has the same sign as either  $f(a_1)$  or  $f(b_1)$ . When  $f(p_1)$  and  $f(a_1)$  have the same sign,  $p \in (p_1, b_1)$ , and we set  $a_2 = p_1$  and  $b_2 = b_1$ . When  $f(p_1)$  and  $f(a_1)$  have opposite signs,  $p \in (a_1, p_1)$ , and we set  $a_2 = a_1$  and  $b_2 = p_1$ . We then reapply the process to the interval  $[a_2, b_2]$ . This produces the method described in Algorithm 2.1. (See Figure 2.1.)

Figure 2.1



ALGORITHM

2.1

**Bisection**

To find a solution to  $f(x) = 0$  given the continuous function  $f$  on the interval  $[a, b]$ , where  $f(a)$  and  $f(b)$  have opposite signs:

INPUT endpoints  $a, b$ ; tolerance  $TOL$ ; maximum number of iterations  $N_0$ .

OUTPUT approximate solution  $p$  or message of failure.

Step 1 Set  $i = 1$ ;  
 $FA = f(a)$ .

Step 2 While  $i \leq N_0$  do Steps 3–6.

Step 3 Set  $p = a + (b - a)/2$ ; (Compute  $p_i$ .)  
 $FP = f(p)$ .

**Step 4** If  $FP = 0$  or  $(b - a)/2 < TOL$  then  
 OUTPUT ( $p$ ); (*Procedure completed successfully.*)  
 STOP.

**Step 5** Set  $i = i + 1$ .

**Step 6** If  $FA \cdot FP > 0$  then set  $a = p$ ; (*Compute  $a_i, b_i$ .*)  
 $FA = FP$   
 else set  $b = p$ .

**Step 7** OUTPUT ('Method failed after  $N_0$  iterations,  $N_0 =$ ',  $N_0$ );  
 (*The procedure was unsuccessful.*)  
 STOP. ■

Other stopping procedures can be applied in Step 4 of Algorithm 2.1 or in any of the iterative techniques in this chapter. For example, we can select a tolerance  $\epsilon > 0$  and generate  $p_1, \dots, p_N$  until one of the following conditions is met:

$$|p_N - p_{N-1}| < \epsilon, \quad (2.1)$$

$$\frac{|p_N - p_{N-1}|}{|p_N|} < \epsilon, \quad p_N \neq 0, \quad \text{or} \quad (2.2)$$

$$|f(p_N)| < \epsilon. \quad (2.3)$$

Unfortunately, difficulties can arise using any of these stopping criteria. For example, there are sequences  $\{p_n\}_{n=0}^{\infty}$  with the property that the differences  $p_n - p_{n-1}$  converge to zero while the sequence itself diverges. (See Exercise 15.) It is also possible for  $f(p_n)$  to be close to zero while  $p_n$  differs significantly from  $p$ . (See Exercise 14.) Without additional knowledge about  $f$  or  $p$ , Inequality (2.2) is the best stopping criterion to apply because it comes closest to testing relative error.

When using a computer to generate approximations, it is good practice to set an upper bound on the number of iterations. This will eliminate the possibility of entering an infinite loop, a situation that can arise when the sequence diverges (and also when the program is incorrectly coded). This was done in Step 2 of Algorithm 2.1 where the bound  $N_0$  was set and the procedure terminated if  $i > N_0$ .

Note that to start the Bisection Algorithm, an interval  $[a, b]$  must be found with  $f(a) \cdot f(b) < 0$ . At each step the length of the interval known to contain a zero of  $f$  is reduced by a factor of 2; hence it is advantageous to choose the initial interval  $[a, b]$  as small as possible. For example, if  $f(x) = 2x^3 - x^2 + x - 1$ , we have both

$$f(-4) \cdot f(4) < 0 \quad \text{and} \quad f(0) \cdot f(1) < 0,$$

so the Bisection Algorithm could be used on either of the intervals  $[-4, 4]$  or  $[0, 1]$ . Starting the Bisection Algorithm on  $[0, 1]$  instead of  $[-4, 4]$  will reduce by 3 the number of iterations required to achieve a specified accuracy.

The following example illustrates the Bisection Algorithm. The iteration in this example is terminated when the relative error is less than 0.0001; that is, when

$$\frac{|p - p_n|}{|p|} < 10^{-4}.$$

**EXAMPLE 1** The equation  $f(x) = x^3 + 4x^2 - 10 = 0$  has a root in  $[1, 2]$  since  $f(1) = -5$  and  $f(2) = 14$ . The Bisection Algorithm gives the values in Table 2.1.

**Table 2.1**

$n$	$a_n$	$b_n$	$p_n$	$f(p_n)$
1	1.0	2.0	1.5	2.375
2	1.0	1.5	1.25	-1.79687
3	1.25	1.5	1.375	0.16211
4	1.25	1.375	1.3125	-0.84839
5	1.3125	1.375	1.34375	-0.35098
6	1.34375	1.375	1.359375	-0.09641
7	1.359375	1.375	1.3671875	0.03236
8	1.359375	1.3671875	1.36328125	-0.03215
9	1.36328125	1.3671875	1.365234375	0.000072
10	1.36328125	1.365234375	1.364257813	-0.01605
11	1.364257813	1.365234375	1.364746094	-0.00799
12	1.364746094	1.365234375	1.364990235	-0.00396
13	1.364990235	1.365234375	1.365112305	-0.00194

After 13 iterations,  $p_{13} = 1.365112305$  approximates the root  $p$  with an error

$$|p - p_{13}| < |b_{14} - a_{14}| = |1.365234375 - 1.365112305| = 0.000122070.$$

Since  $|a_{14}| < |p|$ ,

$$\frac{|p - p_{13}|}{|p|} < \frac{|b_{14} - a_{14}|}{|a_{14}|} \leq 9.0 \times 10^{-5},$$

so the approximation is correct to at least four significant digits. The correct value of  $p$ , to nine decimal places, is  $p = 1.365230013$ . Note that  $p_9$  is closer to  $p$  than is the final approximation  $p_{13}$ . You might suspect this is true since  $|f(p_9)| < |f(p_{13})|$ , but we cannot be sure of this unless the true answer is known. ■

The Bisection method, though conceptually clear, has significant drawbacks. It is slow to converge (that is,  $N$  may become quite large before  $|p - p_N|$  is sufficiently small), and a good intermediate approximation can be inadvertently discarded. However, the method has the important property that it always converges to a solution, and for that reason it is often used as a starter for the more efficient methods we will present later in this chapter.

**Theorem 2.1**

Suppose that  $f \in C[a, b]$  and  $f(a) \cdot f(b) < 0$ . The Bisection method generates a sequence  $\{p_n\}_{n=1}^{\infty}$  approximating a zero  $p$  of  $f$  with

$$|p_n - p| \leq \frac{b - a}{2^n}, \quad \text{when } n \geq 1. \quad \blacksquare$$

**Proof** For each  $n \geq 1$ , we have

$$b_n - a_n = \frac{1}{2^{n-1}}(b - a) \quad \text{and} \quad p \in (a_n, b_n).$$

Since  $p_n = \frac{1}{2}(a_n + b_n)$  for all  $n \geq 1$ , it follows that

$$|p_n - p| \leq \frac{1}{2}(b_n - a_n) = \frac{b - a}{2^n}. \quad \dots$$

Since

$$|p_n - p| \leq (b - a) \frac{1}{2^n},$$

the sequence  $\{p_n\}_{n=1}^{\infty}$  converges to  $p$  with rate of convergence  $O\left(\frac{1}{2^n}\right)$ ; that is,

$$p_n = p + O\left(\frac{1}{2^n}\right).$$

It is important to realize that Theorem 2.1 gives only a bound for approximation error and that this bound may be quite conservative. For example, this bound applied to the problem in Example 1 ensures only that

$$|p - p_9| \leq \frac{2 - 1}{2^9} \approx 2 \times 10^{-3},$$

but the actual error is much smaller:

$$|p - p_9| = |1.365230013 - 1.365234375| \approx 4.4 \times 10^{-6}.$$

**EXAMPLE 2**

To determine the number of iterations necessary to solve  $f(x) = x^3 + 4x^2 - 10 = 0$  with accuracy  $10^{-3}$  using  $a_1 = 1$  and  $b_1 = 2$  requires finding an integer  $N$  that satisfies

$$|p_N - p| \leq 2^{-N}(b - a) = 2^{-N} < 10^{-3}.$$

To determine  $N$  we will use logarithms. Although logarithms to any base would suffice, we will use base-10 logarithms since the tolerance is given as a power of 10. Since  $2^{-N} < 10^{-3}$  implies that  $\log_{10} 2^{-N} < \log_{10} 10^{-3} = -3$ , we have

$$-N \log_{10} 2 < -3 \quad \text{and} \quad N > \frac{3}{\log_{10} 2} \approx 9.96.$$

Hence, ten iterations will ensure an approximation accurate to within  $10^{-3}$ . Table 2.1 on page 51 shows that the value of  $p_9 = 1.365234375$  is accurate to within  $10^{-4}$ . Again, it

is important to keep in mind that the error analysis gives only a bound for the number of iterations, and in many cases this bound is much larger than the actual number required. ■

The bound for the number of iterations for the Bisection method assumes that the calculations are performed using infinite-digit arithmetic. When implementing the method on a computer, consideration must be given to the effects of roundoff error. For example, the computation of the midpoint of the interval  $[a_n, b_n]$  should be found from the equation

$$p_n = a_n + \frac{b_n - a_n}{2}$$

instead of from the algebraically equivalent equation

$$p_n = \frac{a_n + b_n}{2}.$$

The first equation adds a small correction,  $(b_n - a_n)/2$ , to the known value  $a_n$ . When  $b_n - a_n$  is near the maximum precision of the machine this correction might be in error, but the error would not significantly affect the computed value of  $p_n$ . However, when  $b_n - a_n$  is near the maximum precision of the machine, it is possible for  $(a_n + b_n)/2$  to return a midpoint that is not even in the interval  $[a_n, b_n]$ .

As a final remark, to determine which subinterval of  $[a_n, b_n]$  contains a root of  $f$ , it is better to make use of the **signum** function, which is defined as

$$\text{sgn}(x) = \begin{cases} -1, & \text{if } x < 0, \\ 0, & \text{if } x = 0, \\ 1, & \text{if } x > 0. \end{cases}$$

The test

$$\text{sgn}(f(a_n)) \text{sgn}(f(p_n)) < 0 \quad \text{instead of} \quad f(a_n)f(p_n) < 0$$

gives the same result but avoids the possibility of overflow or underflow in the multiplication of  $f(a_n)$  and  $f(p_n)$ .

## EXERCISE SET 2.1

1. Use the Bisection method to find  $p_3$  for  $f(x) = \sqrt{x} - \cos x$  on  $[0, 1]$ .
2. Let  $f(x) = 3(x + 1)(x - \frac{1}{2})(x - 1)$ . Use the Bisection method on the following intervals to find  $p_3$ .
  - a.  $[-2, 1.5]$
  - b.  $[-1.25, 2.5]$
3. Use the Bisection method to find solutions accurate to within  $10^{-2}$  for  $x^3 - 7x^2 + 14x - 6 = 0$  on each interval.
  - a.  $[0, 1]$
  - b.  $[1, 3.2]$
  - c.  $[3.2, 4]$

4. Use the Bisection method to find solutions accurate to within  $10^{-2}$  for  $x^4 - 2x^3 - 4x^2 + 4x + 4 = 0$  on each interval.
  - a.  $[-2, -1]$
  - b.  $[0, 2]$
  - c.  $[2, 3]$
  - d.  $[-1, 0]$
5. Use the Bisection method to find a solution accurate to within  $10^{-3}$  for  $x = \tan x$  on  $[4, 4.5]$ .
6. Use the Bisection method to find a solution accurate to within  $10^{-3}$  for  $2 + \cos(e^x - 2) - e^x = 0$  on  $[0.5, 1.5]$ .
7. Use the Bisection method to find solutions accurate to within  $10^{-5}$  for the following problems.
  - a.  $x - 2^{-x} = 0$  for  $0 \leq x \leq 1$
  - b.  $e^x - x^2 + 3x - 2 = 0$  for  $0 \leq x \leq 1$
  - c.  $2x \cos(2x) - (x + 1)^2 = 0$  for  $-3 \leq x \leq -2$  and  $-1 \leq x \leq 0$
  - d.  $x \cos x - 2x^2 + 3x - 1 = 0$  for  $0.2 \leq x \leq 0.3$  and  $1.2 \leq x \leq 1.3$
8. Let  $f(x) = (x + 2)(x + 1)^2x(x - 1)^3(x - 2)$ . To which zero of  $f$  does the Bisection method converge when applied on the following intervals?
  - a.  $[-1.5, 2.5]$
  - b.  $[-0.5, 2.4]$
  - c.  $[-0.5, 3]$
  - d.  $[-3, -0.5]$
9. Let  $f(x) = (x + 2)(x + 1)x(x - 1)^3(x - 2)$ . To which zero of  $f$  does the Bisection method converge when applied on the following intervals?
  - a.  $[-3, 2.5]$
  - b.  $[-2.5, 3]$
  - c.  $[-1.75, 1.5]$
  - d.  $[-1.5, 1.75]$
10. Find an approximation to  $\sqrt{3}$  correct to within  $10^{-4}$  using the Bisection Algorithm. [Hint: Consider  $f(x) = x^2 - 3$ .]
11. Find an approximation to  $\sqrt[3]{25}$  correct to within  $10^{-4}$  using the Bisection Algorithm.
12. Use Theorem 2.1 to find a bound for the number of iterations needed to achieve an approximation with accuracy  $10^{-3}$  to the solution of  $x^3 + x - 4 = 0$  lying in the interval  $[1, 4]$ . Find an approximation to the root with this degree of accuracy.
13. Use Theorem 2.1 to find a bound for the number of iterations needed to achieve an approximation with accuracy  $10^{-4}$  to the solution of  $x^3 - x - 1 = 0$  lying in the interval  $[1, 2]$ . Find an approximation to the root with this degree of accuracy.
14. Let  $f(x) = (x - 1)^{10}$ ,  $p = 1$ , and  $p_n = 1 + 1/n$ . Show that  $|f(p_n)| < 10^{-3}$  whenever  $n > 1$  but that  $|p - p_n| < 10^{-3}$  requires that  $n > 1000$ .
15. Let  $\{p_n\}$  be the sequence defined by  $p_n = \sum_{k=1}^n \frac{1}{k}$ . Show that  $\{p_n\}$  diverges even though  $\lim_{n \rightarrow \infty} (p_n - p_{n-1}) = 0$ .
16. The function defined by  $f(x) = \sin \pi x$  has zeros at every integer. Show that when  $-1 < a < 0$  and  $2 < b < 3$ , the Bisection method converges to
  - a. 0, if  $a + b < 2$
  - b. 2, if  $a + b > 2$
  - c. 1, if  $a + b = 2$
17. A trough of length  $L$  has a cross section in the shape of a semicircle with radius  $r$  (See the accompanying figure.) When filled with water to within a distance  $h$  of the top, the volume  $V$  of water is

$$V = L [0.5\pi r^2 - r^2 \arcsin(h/r) - h(r^2 - h^2)^{1/2}].$$

