

Solutions of Equations of One Variable

2.1 Introduction

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 $f(x) = 0$. A root of this equation is also called a **zero** of the function f . This is one of the oldest known approximation problems, yet research continues in this area at the present time.

The problem of finding an approximation to the root of an equation can be traced at least as far back 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 given in Exercise 11 of Section 2.4.

2.2 The Bisection Method

The first and most elementary technique we consider is the **Bisection**, or *Binary-Search*, method. The Bisection method is used to determine, to any specified accuracy that your computer will permit, a solution to $f(x) = 0$ on an interval $[a, b]$, provided that f is continuous on the interval and that $f(a)$ and $f(b)$ are of opposite sign. Although the method will work for the case when more than one root is contained in the interval $[a, b]$, we assume for simplicity of our discussion that the root in this interval is unique.

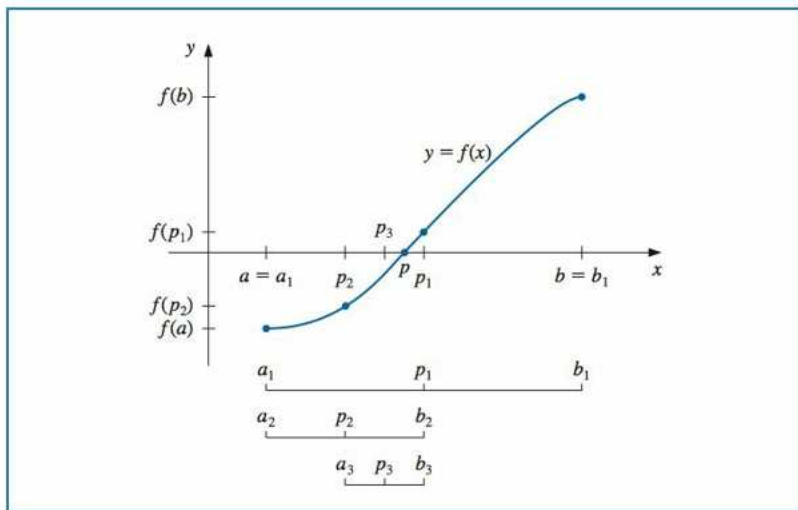
Bisection Technique

To begin the Bisection method, set $a_1 = a$ and $b_1 = b$, as shown in Figure 2.1, and let p_1 be the midpoint of the interval $[a, b]$:

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

If $f(p_1) = 0$, then the root p is given by $p = p_1$; if $f(p_1) \neq 0$, then $f(p_1)$ has the same sign as either $f(a_1)$ or $f(b_1)$.

Figure 2.1



In computer science, the process of dividing a set continually in half to search for the solution to a problem, as the Bisection method does, is known as a *binary search* procedure.

- If $f(p_1)$ and $f(a_1)$ have opposite signs, then p is in the interval (a_1, p_1) , and we set

$$a_2 = a_1 \quad \text{and} \quad b_2 = p_1.$$

- If $f(p_1)$ and $f(a_1)$ have the same sign, then p is in the interval (p_1, b_1) , and we set

$$a_2 = p_1 \quad \text{and} \quad b_2 = b_1.$$

Reapply the process to the interval $[a_2, b_2]$, and continue forming $[a_3, b_3]$, $[a_4, b_4]$, \dots . Each new interval will contain p and have length one half of the length of the preceding interval.

Bisection Method

An interval $[a_{n+1}, b_{n+1}]$ containing an approximation to a root of $f(x) = 0$ is constructed from an interval $[a_n, b_n]$ containing the root by first letting

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

Then set

$$a_{n+1} = a_n \quad \text{and} \quad b_{n+1} = p_n \quad \text{if} \quad f(a_n)f(p_n) < 0,$$

and

$$a_{n+1} = p_n \quad \text{and} \quad b_{n+1} = b_n \quad \text{otherwise.}$$

Program BISECT21
implements the Bisection
method.*

There are three stopping criteria commonly incorporated in the Bisection method, and incorporated within BISECT21.

- The method stops if one of the midpoints happens to coincide with the root.
- It also stops when the length of the search interval is less than some prescribed tolerance we call TOL .
- The procedure also stops if the number of iterations exceeds a preset bound N_0 .

To start the Bisection method, an interval $[a, b]$ must be found with $f(a) \cdot f(b) < 0$; that is, $f(a)$ and $f(b)$ have opposite signs. At each step, the length of the interval known to contain a zero of f is reduced by a factor of 2. Since the midpoint p_1 must be within $(b-a)/2$ of the root p , and each succeeding iteration divides the interval under consideration by 2, we have

$$|p_n - p| \leq \frac{b-a}{2^n}.$$

Consequently, it is easy to determine a bound for the number of iterations needed to ensure a given tolerance. If the root needs to be determined within the tolerance TOL , we need to determine the number of iterations, n , so that

$$\frac{b-a}{2^n} < TOL.$$

Using logarithms to solve for n in this inequality gives

$$\frac{b-a}{TOL} < 2^n, \quad \text{which implies that} \quad \log_2 \left(\frac{b-a}{TOL} \right) < n.$$

Since the number of required iterations to guarantee a given accuracy depends on the length of the initial interval $[a, b]$, we want to choose this interval as small as possible. For example, if $f(x) = 2x^3 - x^2 + x - 1$, we have both

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

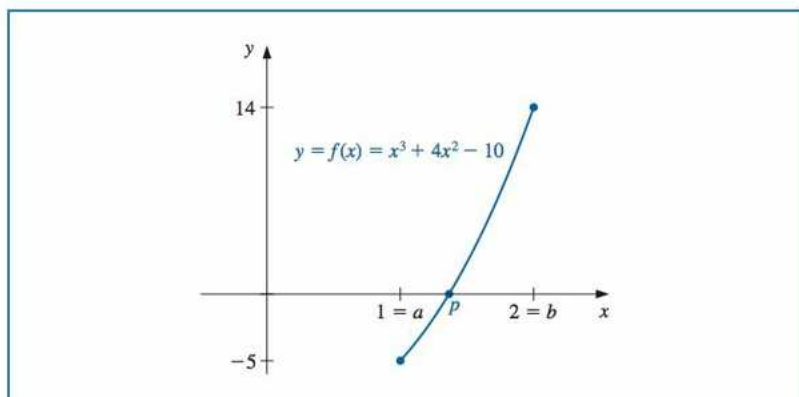
so the Bisection method could be used on either $[-4, 4]$ or $[0, 1]$. Starting the Bisection method on $[0, 1]$ instead of $[-4, 4]$ reduces by 3 the number of iterations required to achieve a specified accuracy.

Example 1 Show that $f(x) = x^3 + 4x^2 - 10 = 0$ has a root in $[1, 2]$ and use the Bisection method to determine an approximation to the root that is accurate to at least within 10^{-4} .

Solution Because $f(1) = -5$ and $f(2) = 14$, the Intermediate Value Theorem ensures that this continuous function has a root in $[1, 2]$. Since $f'(x) = 3x^2 + 8x$ is always positive on $[1, 2]$, the function f is increasing, and, as seen in Figure 2.2, the root is unique.

*These programs can be found at <http://www.math.ysu.edu/~fares/Numerical-Methods/Programs/>

Figure 2.2



For the first iteration of the Bisection method we use the fact that at the midpoint of $[1, 2]$ we have

$$f(1.5) = 2.375 > 0.$$

This indicates that we should select the interval $[1, 1.5]$ for our second iteration. Then we find that

$$f(1.25) = -1.796875$$

so our new interval becomes $[1.25, 1.5]$, whose midpoint is 1.375. Continuing in this manner 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 $|a_{14}| < |p|$, we have

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

so the approximation is correct to at least within 10^{-4} . 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 because $|f(p_9)| < |f(p_{13})|$. ■

The Bisection method, although conceptually clear, has serious drawbacks. It is slow to converge relative to the other techniques we will discuss, and a good intermediate approximation may be inadvertently discarded. This happened, for example, with p_9 in Example 1. However, the method has the important property that it always converges to a solution and it is easy to determine a bound for the number of iterations needed to ensure a given accuracy. For these reasons, the Bisection method is frequently used as a dependable starting procedure for the more efficient methods presented later in this chapter.

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 round-off 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, in the second equation, if $b_n - a_n$ is near the maximum precision of the machine, it is possible for p_n to return a midpoint that is not even in the interval $[a_n, b_n]$.

A number of tests can be used to see if a root has been found. We would normally use a test of the form

$$|f(p_n)| < \epsilon,$$

where $\epsilon > 0$ would be a small number related in some way to the tolerance. However, it is also possible for the value $f(p_n)$ to be small when p_n is not near the root p .

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 **sgn** 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 $f(p_n)$.

The Latin word *signum* means “token” or “sign.” So the **sgn** function quite naturally returns the sign of a number (unless the number is 0).

EXERCISE SET 2.2

- Use the Bisection method to find p_3 for $f(x) = \sqrt{x} - \cos x$ on $[0, 1]$.
- Let $f(x) = 3(x+1)(x-\frac{1}{2})(x-1)$. Use the Bisection method on the following intervals to find p_3 .
 - $[-2, 1.5]$
 - $[-1.25, 2.5]$
- Use the Bisection method to find solutions accurate to within 10^{-2} for $x^3 - 7x^2 + 14x - 6 = 0$ on each interval.
 - $[0, 1]$
 - $[1, 3.2]$
 - $[3.2, 4]$
- 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.
 - $[-2, -1]$
 - $[0, 2]$
 - $[2, 3]$
 - $[-1, 0]$
- Sketch the graphs of $y = x$ and $y = 2 \sin x$.
 - Use the Bisection method to find an approximation to within 10^{-2} to the first positive value of x with $x = 2 \sin x$.
- Sketch the graphs of $y = x$ and $y = \tan x$.
 - Use the Bisection method to find an approximation to within 10^{-2} to the first positive value of x with $x = \tan x$.
- Let $f(x) = (x+2)(x+1)x(x-1)^3(x-2)$. To which zero of f does the Bisection method converge for the following intervals?
 - $[-3, 2.5]$
 - $[-2.5, 3]$
 - $[-1.75, 1.5]$
 - $[-1.5, 1.75]$
- Let $f(x) = (x+2)(x+1)^2x(x-1)^3(x-2)$. To which zero of f does the Bisection method converge for the following intervals?
 - $[-1.5, 2.5]$
 - $[-0.5, 2.4]$
 - $[-0.5, 3]$
 - $[-3, -0.5]$
- Use the Bisection method to find an approximation to $\sqrt{3}$ correct to within 10^{-4} . [Hint: Consider $f(x) = x^2 - 3$.]
- Use the Bisection method to find an approximation to $\sqrt[3]{25}$ correct to within 10^{-4} .
- Find a bound for the number of Bisection method 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.
- Find a bound for the number of Bisection method 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.
- 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
 - 0, if $a + b < 2$
 - 2, if $a + b > 2$
 - 1, if $a + b = 2$

2.3 The Secant Method

Although the Bisection method always converges, the speed of convergence is often too slow for general use. Figure 2.3 gives a graphical interpretation of the Bisection method that can be used to discover how improvements on this technique can be derived. It shows the graph of a continuous function that is negative at a_1 and positive at b_1 . The first approximation p_1 to the root p is found by drawing the line joining the points $(a_1, \text{sgn}(f(a_1))) = (a_1, -1)$ and $(b_1, \text{sgn}(f(b_1))) = (b_1, 1)$ and letting p_1 be the point where this line intersects the x -axis. In essence, the line joining $(a_1, -1)$ and $(b_1, 1)$ has been used to approximate the graph of f on the interval $[a_1, b_1]$. Successive approximations apply this same process on subintervals