

$$q_1 = \int_0^{0.1} f(x) dx = 0.099833416647 \quad \text{and} \quad q_2 = \int_0^{0.1} P_3(x) dx = 0.0998333333333,$$

with error  $0.83314 \times 10^{-7} = 8.3314 \times 10^{-8}$ .

Parts (a) and (b) of the example show how two techniques can produce the same approximation but have differing accuracy assurances. Remember that determining approximations is only part of our objective. An equally important part is to determine a bound for the error of the approximation.

## RCISE SET 1.1

- Show that the following equations have at least one solution in the given intervals.
  - $x \cos x - 2x^2 + 3x - 1 = 0$ ,  $[0.2, 0.3]$  and  $[1.2, 1.3]$
  - $(x - 2)^2 - \ln x = 0$ ,  $[1, 2]$  and  $[e, 4]$
  - $2x \cos(2x) - (x - 2)^2 = 0$ ,  $[2, 3]$  and  $[3, 4]$
  - $x - (\ln x)^x = 0$ ,  $[4, 5]$
- Find intervals containing solutions to the following equations.
  - $x - 3^{-x} = 0$
  - $4x^2 - e^x = 0$
  - $x^3 - 2x^2 - 4x + 2 = 0$
  - $x^3 + 4.001x^2 + 4.002x + 1.101 = 0$
- Show that  $f'(x)$  is 0 at least once in the given intervals.
  - $f(x) = 1 - e^x + (e - 1) \sin((\pi/2)x)$ ,  $[0, 1]$
  - $f(x) = (x - 1) \tan x + x \sin \pi x$ ,  $[0, 1]$
  - $f(x) = x \sin \pi x - (x - 2) \ln x$ ,  $[1, 2]$
  - $f(x) = (x - 2) \sin x \ln(x + 2)$ ,  $[-1, 3]$
- Find  $\max_{a \leq x \leq b} |f(x)|$  for the following functions and intervals.
  - $f(x) = (2 - e^x + 2x)/3$ ,  $[0, 1]$
  - $f(x) = (4x - 3)/(x^2 - 2x)$ ,  $[0.5, 1]$
  - $f(x) = 2x \cos(2x) - (x - 2)^2$ ,  $[2, 4]$
  - $f(x) = 1 + e^{-\cos(x-1)}$ ,  $[1, 2]$
- Use the Intermediate Value Theorem and Rolle's Theorem to show that the graph of  $f(x) = x^3 + 2x + k$  crosses the  $x$ -axis exactly once, regardless of the value of the constant  $k$ .
- Suppose  $f \in C[a, b]$  and  $f'(x)$  exists on  $(a, b)$ . Show that if  $f'(x) \neq 0$  for all  $x$  in  $(a, b)$ , then there can exist at most one number  $p$  in  $[a, b]$  with  $f(p) = 0$ .
- Let  $f(x) = x^3$ .
  - Find the second Taylor polynomial  $P_2(x)$  about  $x_0 = 0$ .
  - Find  $R_2(0.5)$  and the actual error in using  $P_2(0.5)$  to approximate  $f(0.5)$ .
  - Repeat part (a) using  $x_0 = 1$ .
  - Repeat part (b) using the polynomial from part (c).
- Find the third Taylor polynomial  $P_3(x)$  for the function  $f(x) = \sqrt{x + 1}$  about  $x_0 = 0$ . Approximate  $\sqrt{0.5}$ ,  $\sqrt{0.75}$ ,  $\sqrt{1.25}$ , and  $\sqrt{1.5}$  using  $P_3(x)$ , and find the actual errors.
- Find the second Taylor polynomial  $P_2(x)$  for the function  $f(x) = e^x \cos x$  about  $x_0 = 0$ .
  - Use  $P_2(0.5)$  to approximate  $f(0.5)$ . Find an upper bound for error  $|f(0.5) - P_2(0.5)|$  using the error formula, and compare it to the actual error.

- b. Find a bound for the error  $|f(x) - P_2(x)|$  in using  $P_2(x)$  to approximate  $f(x)$  on the interval  $[0, 1]$ .
- c. Approximate  $\int_0^1 f(x) dx$  using  $\int_0^1 P_2(x) dx$ .
- d. Find an upper bound for the error in (c) using  $\int_0^1 |R_2(x) dx|$ , and compare the bound to the actual error.
10. Repeat Exercise 9 using  $x_0 = \pi/6$ .
11. Find the third Taylor polynomial  $P_3(x)$  for the function  $f(x) = (x - 1) \ln x$  about  $x_0 = 1$ .
- a. Use  $P_3(0.5)$  to approximate  $f(0.5)$ . Find an upper bound for error  $|f(0.5) - P_3(0.5)|$  using the error formula, and compare it to the actual error.
- b. Find a bound for the error  $|f(x) - P_3(x)|$  in using  $P_3(x)$  to approximate  $f(x)$  on the interval  $[0.5, 1.5]$ .
- c. Approximate  $\int_{0.5}^{1.5} f(x) dx$  using  $\int_{0.5}^{1.5} P_3(x) dx$ .
- d. Find an upper bound for the error in (c) using  $\int_{0.5}^{1.5} |R_3(x) dx|$ , and compare the bound to the actual error.
12. Let  $f(x) = 2x \cos(2x) - (x - 2)^2$  and  $x_0 = 0$ .
- a. Find the third Taylor polynomial  $P_3(x)$ , and use it to approximate  $f(0.4)$ .
- b. Use the error formula in Taylor's Theorem to find an upper bound for the error  $|f(0.4) - P_3(0.4)|$ . Compute the actual error.
- c. Find the fourth Taylor polynomial  $P_4(x)$ , and use it to approximate  $f(0.4)$ .
- d. Use the error formula in Taylor's Theorem to find an upper bound for the error  $|f(0.4) - P_4(0.4)|$ . Compute the actual error.
13. Find the fourth Taylor polynomial  $P_4(x)$  for the function  $f(x) = xe^{x^2}$  about  $x_0 = 0$ .
- a. Find an upper bound for  $|f(x) - P_4(x)|$ , for  $0 \leq x \leq 0.4$ .
- b. Approximate  $\int_0^{0.4} f(x) dx$  using  $\int_0^{0.4} P_4(x) dx$ .
- c. Find an upper bound for the error in (b) using  $\int_0^{0.4} P_4(x) dx$ .
- d. Approximate  $f'(0.2)$  using  $P_4'(0.2)$ , and find the error.
14. Use the error term of a Taylor polynomial to estimate the error involved in using  $\sin x \approx x$  to approximate  $\sin 1^\circ$ .
15. Use a Taylor polynomial about  $\pi/4$  to approximate  $\cos 42^\circ$  to an accuracy of  $10^{-6}$ .
16. Let  $f(x) = e^{x/2} \sin(x/3)$ . Use Maple to determine the following.
- a. The third Maclaurin polynomial  $P_3(x)$ .
- b.  $f^{(4)}(x)$  and a bound for the error  $|f(x) - P_3(x)|$  on  $[0, 1]$ .
17. Let  $f(x) = \ln(x^2 + 2)$ . Use Maple to determine the following.
- a. The Taylor polynomial  $P_3(x)$  for  $f$  expanded about  $x_0 = 1$ .
- b. The maximum error  $|f(x) - P_3(x)|$ , for  $0 \leq x \leq 1$ .
- c. The Maclaurin polynomial  $\tilde{P}_3(x)$  for  $f$ .
- d. The maximum error  $|f(x) - \tilde{P}_3(x)|$ , for  $0 \leq x \leq 1$ .
- e. Does  $P_3(0)$  approximate  $f(0)$  better than  $\tilde{P}_3(1)$  approximates  $f(1)$ ?
18. Let  $f(x) = (1 - x)^{-1}$  and  $x_0 = 0$ . Find the  $n$ th Taylor polynomial  $P_n(x)$  for  $f(x)$  about  $x_0$ . Find a value of  $n$  necessary for  $P_n(x)$  to approximate  $f(x)$  to within  $10^{-6}$  on  $[0, 0.5]$ .
19. Let  $f(x) = e^x$  and  $x_0 = 0$ . Find the  $n$ th Taylor polynomial  $P_n(x)$  for  $f(x)$  about  $x_0$ . Find a value of  $n$  necessary for  $P_n(x)$  to approximate  $f(x)$  to within  $10^{-6}$  on  $[0, 0.5]$ .
20. Find the  $n$ th Maclaurin polynomial  $P_n(x)$  for  $f(x) = \arctan x$ .
21. The polynomial  $P_2(x) = 1 - \frac{1}{2}x^2$  is to be used to approximate  $f(x) = \cos x$  in  $[-\frac{1}{2}, \frac{1}{2}]$ . Find a bound for the maximum error.

22. The  $n$ th Taylor polynomial for a function  $f$  at  $x_0$  is sometimes referred to as the polynomial of degree at most  $n$  that “best” approximates  $f$  near  $x_0$ .
- Explain why this description is accurate.
  - Find the quadratic polynomial that best approximates a function  $f$  near  $x_0 = 1$  if the tangent line at  $x_0 = 1$  has equation  $y = 4x - 1$ , and if  $f''(1) = 6$ .
23. A Maclaurin polynomial for  $e^x$  is used to give the approximation 2.5 to  $e$ . The error bound in this approximation is established to be  $E = \frac{1}{6}$ . Find a bound for the error in  $E$ .
24. The *error function* defined by

$$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$$

gives the probability that any one of a series of trials will lie within  $x$  units of the mean, assuming that the trials have a normal distribution with mean 0 and standard deviation  $\sqrt{2}/2$ . This integral cannot be evaluated in terms of elementary functions, so an approximating technique must be used.

- Integrate the Maclaurin series for  $e^{-x^2}$  to show that

$$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \sum_{k=0}^{\infty} \frac{(-1)^k x^{2k+1}}{(2k+1)k!}$$

- The error function can also be expressed in the form

$$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} e^{-x^2} \sum_{k=0}^{\infty} \frac{2^k x^{2k+1}}{1 \cdot 3 \cdot 5 \cdots (2k+1)}$$

Verify that the two series agree for  $k = 1, 2, 3$ , and 4. [Hint: Use the Maclaurin series for  $e^{-x^2}$ .]

- Use the series in part (a) to approximate  $\operatorname{erf}(1)$  to within  $10^{-7}$ .
  - Use the same number of terms as in part (c) to approximate  $\operatorname{erf}(1)$  with the series in part (b).
  - Explain why difficulties occur using the series in part (b) to approximate  $\operatorname{erf}(x)$ .
25. A function  $f : [a, b] \rightarrow \mathbb{R}$  is said to satisfy a *Lipschitz condition* with Lipschitz constant  $L$  on  $[a, b]$  if, for every  $x, y \in [a, b]$ , we have  $|f(x) - f(y)| \leq L|x - y|$ .
- Show that if  $f$  satisfies a Lipschitz condition with Lipschitz constant  $L$  on an interval  $[a, b]$ , then  $f \in C[a, b]$ .
  - Show that if  $f$  has a derivative that is bounded on  $[a, b]$  by  $L$ , then  $f$  satisfies a Lipschitz condition with Lipschitz constant  $L$  on  $[a, b]$ .
  - Give an example of a function that is continuous on a closed interval but does not satisfy a Lipschitz condition on the interval.
26. Suppose  $f \in C[a, b]$ , that  $x_1$  and  $x_2$  are in  $[a, b]$ , and that  $c_1$  and  $c_2$  are positive constants. Show that a number  $\xi$  exists between  $x_1$  and  $x_2$  with

$$f(\xi) = \frac{c_1 f(x_1) + c_2 f(x_2)}{c_1 + c_2}$$

27. Let  $f \in C[a, b]$ , and let  $p$  be in the open interval  $(a, b)$ .
- Suppose  $f(p) \neq 0$ . Show that a  $\delta > 0$  exists with  $f(x) \neq 0$  for all  $x$  in  $[p - \delta, p + \delta]$ , where  $[p - \delta, p + \delta]$  a subset of  $[a, b]$ .
  - Suppose  $f(p) = 0$  and  $k > 0$  is given. Show that a  $\delta > 0$  exists with  $|f(x)| \leq k$  for all  $x$  in  $[p - \delta, p + \delta]$ , where  $[p - \delta, p + \delta]$  a subset of  $[a, b]$ .

Nesting has reduced the relative error for the chopping approximation to less than 10% of that obtained initially. For the rounding approximation, the improvement has been even more dramatic; the error in this case has been reduced by more than 95%. ■

Polynomials should *always* be expressed in nested form before performing an evaluation, because this form minimizes the number of arithmetic calculations. The decreased error in Example 6 is due to the reduction in computations from four multiplications and three additions to two multiplications and three additions. One way to reduce round-off error is to reduce the number of error-producing computations.

## XERCISE SET 1.2

- Compute the absolute error and relative error in approximations of  $p$  by  $p^*$ .
  - $p = \pi$ ,  $p^* = 22/7$
  - $p = \pi$ ,  $p^* = 3.1416$
  - $p = e$ ,  $p^* = 2.718$
  - $p = \sqrt{2}$ ,  $p^* = 1.414$
  - $p = e^{10}$ ,  $p^* = 22000$
  - $p = 10^\pi$ ,  $p^* = 1400$
  - $p = 8!$ ,  $p^* = 39900$
  - $p = 9!$ ,  $p^* = \sqrt{18\pi}(9/e)^9$
- Find the largest interval in which  $p^*$  must lie to approximate  $p$  with relative error at most  $10^{-4}$  for each value of  $p$ .
  - $\pi$
  - $e$
  - $\sqrt{2}$
  - $\sqrt[3]{7}$
- Suppose  $p^*$  must approximate  $p$  with relative error at most  $10^{-3}$ . Find the largest interval in which  $p^*$  must lie for each value of  $p$ .
  - 150
  - 900
  - 1500
  - 90
- Perform the following computations (i) exactly, (ii) using three-digit chopping arithmetic, and (iii) using three-digit rounding arithmetic. (iv) Compute the relative errors in parts (ii) and (iii).
  - $\frac{4}{5} + \frac{1}{3}$
  - $\frac{4}{5} \cdot \frac{1}{3}$
  - $(\frac{1}{3} - \frac{3}{11}) + \frac{3}{20}$
  - $(\frac{1}{3} + \frac{3}{11}) - \frac{3}{20}$
- Use three-digit rounding arithmetic to perform the following calculations. Compute the absolute error and relative error with the exact value determined to at least five digits.
  - $133 + 0.921$
  - $133 - 0.499$
  - $(121 - 0.327) - 119$
  - $(121 - 119) - 0.327$
  - $\frac{\frac{13}{14} - \frac{6}{7}}{2e - 5.4}$
  - $-10\pi + 6e - \frac{3}{62}$
  - $(\frac{2}{9}) \cdot (\frac{9}{7})$
  - $\frac{\pi - \frac{22}{7}}{\frac{1}{17}}$
- Repeat Exercise 5 using four-digit rounding arithmetic.
- Repeat Exercise 5 using three-digit chopping arithmetic.
- Repeat Exercise 5 using four-digit chopping arithmetic.
- The first three nonzero terms of the Maclaurin series for the arctangent function are  $x - (1/3)x^3 + (1/5)x^5$ . Compute the absolute error and relative error in the following approximations of  $\pi$  using the polynomial in place of the arctangent:
  - $4[\arctan(\frac{1}{2}) + \arctan(\frac{1}{3})]$
  - $16\arctan(\frac{1}{5}) - 4\arctan(\frac{1}{239})$
- The number  $e$  can be defined by  $e = \sum_{n=0}^{\infty} (1/n!)$ . Compute the absolute error and relative error in the following approximations of  $e$ :

a.  $\sum_{n=0}^5 \frac{1}{n!}$

b.  $\sum_{n=0}^{10} \frac{1}{n!}$

11. Let

$$f(x) = \frac{x \cos x - \sin x}{x - \sin x}.$$

- Find  $\lim_{x \rightarrow 0} f(x)$ .
- Use four-digit rounding arithmetic to evaluate  $f(0.1)$ .
- Replace each trigonometric function with its third Maclaurin polynomial, and repeat part (b).
- The actual value is  $f(0.1) = -1.99899998$ . Find the relative error for the values obtained in parts (b) and (c).

12. Let

$$f(x) = \frac{e^x - e^{-x}}{x}.$$

- Find  $\lim_{x \rightarrow 0} (e^x - e^{-x})/x$ .
  - Use three-digit rounding arithmetic to evaluate  $f(0.1)$ .
  - Replace each exponential function with its third Maclaurin polynomial, and repeat part (b).
  - The actual value is  $f(0.1) = 2.003335000$ . Find the relative error for the values obtained in parts (b) and (c).
13. Use four-digit rounding arithmetic and the formulas of Example 5 to find the most accurate approximations to the roots of the following quadratic equations. Compute the absolute errors and relative errors.
- $\frac{1}{3}x^2 - \frac{123}{4}x + \frac{1}{6} = 0$
  - $\frac{1}{3}x^2 + \frac{123}{4}x - \frac{1}{6} = 0$
  - $1.002x^2 - 11.01x + 0.01265 = 0$
  - $1.002x^2 + 11.01x + 0.01265 = 0$
14. Repeat Exercise 13 using four-digit chopping arithmetic.
15. Use the 64-bit long real format to find the decimal equivalent of the following floating-point machine numbers.
- 0 10000001010 10010011000000000000000000 000000000000000000000000000000
  - 1 10000001010 10010011000000000000000000 000000000000000000000000000000
  - 0 01111111111 01010011000000000000000000 000000000000000000000000000000
  - 0 01111111111 01010011000000000000000000 000000000000000000000000000001
16. Find the next largest and smallest machine numbers in decimal form for the numbers given in Exercise 15.
17. Suppose two points  $(x_0, y_0)$  and  $(x_1, y_1)$  are on a straight line with  $y_1 \neq y_0$ . Two formulas are available to find the  $x$ -intercept of the line:

$$x = \frac{x_0 y_1 - x_1 y_0}{y_1 - y_0} \quad \text{and} \quad x = x_0 - \frac{(x_1 - x_0)y_0}{y_1 - y_0}.$$

- Show that both formulas are algebraically correct.
- Use the data  $(x_0, y_0) = (1.31, 3.24)$  and  $(x_1, y_1) = (1.93, 4.76)$  and three-digit rounding arithmetic to compute the  $x$ -intercept both ways. Which method is better and why?

18. The Taylor polynomial of degree  $n$  for  $f(x) = e^x$  is  $\sum_{i=0}^n (x^i/i!)$ . Use the Taylor polynomial of degree nine and three-digit chopping arithmetic to find an approximation to  $e^{-5}$  by each of the following methods.

a. 
$$e^{-5} \approx \sum_{i=0}^9 \frac{(-5)^i}{i!} = \sum_{i=0}^9 \frac{(-1)^i 5^i}{i!}$$

b. 
$$e^{-5} = \frac{1}{e^5} \approx \frac{1}{\sum_{i=0}^9 \frac{5^i}{i!}}$$

- c. An approximate value of  $e^{-5}$  correct to three digits is  $6.74 \times 10^{-3}$ . Which formula, (a) or (b), gives the most accuracy, and why?
19. The two-by-two linear system

$$ax + by = e, \quad cx + dy = f,$$

where  $a, b, c, d, e, f$  are given, can be solved for  $x$  and  $y$  as follows:

$$\text{set } m = \frac{c}{a}, \quad \text{provided } a \neq 0;$$

$$d_1 = d - mb;$$

$$f_1 = f - me;$$

$$y = \frac{f_1}{d_1};$$

$$x = \frac{e - by}{a}.$$

Solve the following linear systems using four-digit rounding arithmetic.

a.  $1.130x - 6.990y = 14.20$

b.  $8.110x + 12.20y = -0.1370$

$1.013x - 6.099y = 14.22$

$-18.11x + 112.2y = -0.1376$

20. Repeat Exercise 19 using four-digit chopping arithmetic.
21. a. Show that the polynomial nesting technique described in Example 6 can also be applied to the evaluation of

$$f(x) = 1.01e^{4x} - 4.62e^{3x} - 3.11e^{2x} + 12.2e^x - 1.99.$$

- b. Use three-digit rounding arithmetic, the assumption that  $e^{1.53} = 4.62$ , and the fact that  $e^{nx} = (e^x)^n$  to evaluate  $f(1.53)$  as given in part (a).
- c. Redo the calculation in part (b) by first nesting the calculations.
- d. Compare the approximations in parts (b) and (c) to the true three-digit result  $f(1.53) = -7.61$ .
22. A rectangular parallelepiped has sides of length 3 cm, 4 cm, and 5 cm, measured to the nearest centimeter. What are the best upper and lower bounds for the volume of this parallelepiped? What are the best upper and lower bounds for the surface area?
23. Let  $P_n(x)$  be the Maclaurin polynomial of degree  $n$  for the arctangent function. Use Maple carrying 75 decimal digits to find the value of  $n$  required to approximate  $\pi$  to within  $10^{-25}$ , using the following formulas.

a.  $4\left[P_n\left(\frac{1}{2}\right) + P_n\left(\frac{1}{3}\right)\right]$

b.  $16P_n\left(\frac{1}{5}\right) - 4P_n\left(\frac{1}{239}\right)$

24. Suppose that  $fl(y)$  is a  $k$ -digit rounding approximation to  $y$ . Show that

$$\left| \frac{y - fl(y)}{y} \right| \leq 0.5 \times 10^{-k+1}.$$

[Hint: If  $d_{k+1} < 5$ , then  $fl(y) = 0.d_1d_2 \dots d_k \times 10^n$ . If  $d_{k+1} \geq 5$ , then  $fl(y) = 0.d_1d_2 \dots d_k \times 10^n + 10^{n-k}$ .]

25. The binomial coefficient

$$\binom{m}{k} = \frac{m!}{k!(m-k)!}$$

describes the number of ways of choosing a subset of  $k$  objects from a set of  $m$  elements.

- a. Suppose decimal machine numbers are of the form

$$\pm 0.d_1d_2d_3d_4 \times 10^n, \quad \text{with } 1 \leq d_1 \leq 9, 0 \leq d_i \leq 9, \text{ if } i = 2, 3, 4 \quad \text{and} \quad |n| \leq 15.$$

What is the largest value of  $m$  for which the binomial coefficient  $\binom{m}{k}$  can be computed for all  $k$  by the definition without causing overflow?

- b. Show that  $\binom{m}{k}$  can also be computed by

$$\binom{m}{k} = \binom{m}{k} \binom{m-1}{k-1} \dots \binom{m-k+1}{1}.$$

- c. What is the largest value of  $m$  for which the binomial coefficient  $\binom{m}{3}$  can be computed by the formula in part (b) without causing overflow?
- d. Use the equation in (b) and four-digit chopping arithmetic to compute the number of possible 5-card hands in a 52-card deck. Compute the actual and relative errors.

26. Let  $f \in C[a, b]$  be a function whose derivative exists on  $(a, b)$ . Suppose  $f$  is to be evaluated at  $x_0$  in  $(a, b)$ , but instead of computing the actual value  $f(x_0)$ , the approximate value,  $\tilde{f}(x_0)$ , is the actual value of  $f$  at  $x_0 + \epsilon$ , that is,  $\tilde{f}(x_0) = f(x_0 + \epsilon)$ .

- a. Use the Mean Value Theorem to estimate the absolute error  $|f(x_0) - \tilde{f}(x_0)|$  and the relative error  $|f(x_0) - \tilde{f}(x_0)|/|f(x_0)|$ , assuming  $f(x_0) \neq 0$ .
- b. If  $\epsilon = 5 \times 10^{-6}$  and  $x_0 = 1$ , find bounds for the absolute and relative errors for
- $f(x) = e^x$
  - $f(x) = \sin x$
- c. Repeat part (b) with  $\epsilon = (5 \times 10^{-6})x_0$  and  $x_0 = 10$ .

27. The following Maple procedure chops a floating-point number  $x$  to  $t$  digits.

```
chop:=proc(x,t);
  if x=0 then 0
  else
    e:=ceil(evalf(log10(abs(x))));
    x2:=evalf(trunc(x*10^(t-e))
    *10^(e-t));
  fi
end;
```

Verify the procedure works for the following values.

- |                          |                          |
|--------------------------|--------------------------|
| a. $x = 124.031, t = 5$  | b. $x = 124.036, t = 5$  |
| c. $x = -124.031, t = 5$ | d. $x = -124.036, t = 5$ |
| e. $x = 0.00653, t = 2$  | f. $x = 0.00656, t = 2$  |
| g. $x = -0.00653, t = 2$ | h. $x = -0.00656, t = 2$ |

28. The opening example to this chapter described a physical experiment involving the temperature of a gas under pressure. In this application, we were given  $P = 1.00$  atm,  $V = 0.100$  m<sup>3</sup>,  $N = 0.00420$

mol, and  $R = 0.08206$ . Solving for  $T$  in the ideal gas law gives

$$T = \frac{PV}{NR} = \frac{(1.00)(0.100)}{(0.00420)(0.08206)} = 290.15 \text{ K} = 17^\circ\text{C}.$$

In the laboratory, it was found that  $T$  was  $15^\circ\text{C}$  under these conditions, and when the pressure was doubled and the volume halved,  $T$  was  $19^\circ\text{C}$ . Assume that the data are rounded values accurate to the places given, and show that both laboratory figures are within the bounds of accuracy for the ideal gas law.

## 1.3 Algorithms and Convergence

Throughout the text we will examine approximation procedures, called *algorithms*, involving sequences of calculations. An **algorithm** is a procedure that describes, in an unambiguous manner, a finite sequence of steps to be performed in a specified order. The object of the algorithm is to implement a procedure to solve a problem or approximate a solution to the problem.

We use a **pseudocode** to describe the algorithms. This pseudocode specifies the form of the input to be supplied and the form of the desired output. Not all numerical procedures give satisfactory output for arbitrarily chosen input. As a consequence, a stopping technique independent of the numerical technique is incorporated into each algorithm to avoid infinite loops.

Two punctuation symbols are used in the algorithms:

A period (.) indicates the termination of a step,

a semicolon (;) separates tasks within a step.

Indentation is used to indicate that groups of statements are to be treated as a single entity.

Looping techniques in the algorithms are either counter-controlled, such as

For  $i = 1, 2, \dots, n$

Set  $x_i = a_i + i \cdot h$

or condition-controlled, such as

While  $i < N$  do Steps 3–6.

To allow for conditional execution, we use the standard

If...then      or      If...      then  
else

constructions.

The steps in the algorithms follow the rules of structured program construction. They have been arranged so that there should be minimal difficulty translating pseudocode into any programming language suitable for scientific applications.

The algorithms are liberally laced with comments. These are written in italics and contained within parentheses to distinguish them from the algorithmic statements.

The use of an algorithm is as old as formal mathematics, but the name derives from the Arabic mathematician Muhammad ibn-Muṣā al-Khwārizmī (c. 780–850). The Latin translation of his works begins with the words "Algoritmi" meaning "Al-Khwārizmī says."



We also use big oh notation to describe the rate at which functions converge.

**Definition 1.19** Suppose that  $\lim_{h \rightarrow 0} G(h) = 0$  and  $\lim_{h \rightarrow 0} F(h) = L$ . If a positive constant  $K$  exists with

$$|F(h) - L| \leq K|G(h)|, \quad \text{for sufficiently small } h,$$

then we write  $F(h) = L + O(G(h))$ . ■

The functions we use for comparison generally have the form  $G(h) = h^p$ , where  $p > 0$ . We are interested in the largest value of  $p$  for which  $F(h) = L + O(h^p)$ .

**Example 5** In Example 3(b) of Section 1.1 we found that the third Taylor polynomial gives

$$\cos h = 1 - \frac{1}{2}h^2 + \frac{1}{24}h^4 \cos \tilde{\xi}(h),$$

for some number  $\tilde{\xi}(h)$  between zero and  $h$ . Consequently,

$$\cos h + \frac{1}{2}h^2 = 1 + \frac{1}{24}h^4 \cos \tilde{\xi}(h).$$

This implies that

$$\cos h + \frac{1}{2}h^2 = 1 + O(h^4),$$

since  $|\cos h + \frac{1}{2}h^2 - 1| = |\frac{1}{24} \cos \tilde{\xi}(h)|h^4 \leq \frac{1}{24}h^4$ . The implication is that as  $h \rightarrow 0$ ,  $\cos h + \frac{1}{2}h^2$  converges to its limit, 1, about as fast as  $h^4$  converges to 0. ■

## EXERCISE SET 1.3

- Use three-digit chopping arithmetic to compute the sum  $\sum_{i=1}^{10} (1/i^2)$  first by  $\frac{1}{1} + \frac{1}{4} + \dots + \frac{1}{100}$  and then by  $\frac{1}{100} + \frac{1}{81} + \dots + \frac{1}{1}$ . Which method is more accurate, and why?
  - Write an algorithm to sum the finite series  $\sum_{i=1}^N x_i$  in reverse order.
- The number  $e$  is defined by  $e = \sum_{n=0}^{\infty} (1/n!)$ . Use four-digit chopping arithmetic to compute the following approximations to  $e$ , and determine the absolute and relative errors.

a.  $e \approx \sum_{n=0}^5 \frac{1}{n!}$

b.  $e \approx \sum_{j=0}^5 \frac{1}{(5-j)!}$

c.  $e \approx \sum_{n=0}^{10} \frac{1}{n!}$

d.  $e \approx \sum_{j=0}^{10} \frac{1}{(10-j)!}$

- The Maclaurin series for the arctangent function converges for  $-1 < x \leq 1$  and is given by

$$\arctan x = \lim_{n \rightarrow \infty} P_n(x) = \lim_{n \rightarrow \infty} \sum_{i=1}^n (-1)^{i+1} \frac{x^{2i-1}}{2i-1}.$$

- Use the fact that  $\tan \pi/4 = 1$  to determine the number of  $n$  terms of the series that need to be summed to ensure that  $|4P_n(1) - \pi| < 10^{-3}$ .
- The C++ programming language requires the value of  $\pi$  to be within  $10^{-10}$ . How many terms of the series would we need to sum to obtain this degree of accuracy?

4. Exercise 3 details a rather inefficient means of obtaining an approximation to  $\pi$ . The method can be improved substantially by observing that  $\pi/4 = \arctan \frac{1}{2} + \arctan \frac{1}{3}$  and evaluating the series for the arctangent at  $\frac{1}{2}$  and at  $\frac{1}{3}$ . Determine the number of terms that must be summed to ensure an approximation to  $\pi$  to within  $10^{-3}$ .
5. Another formula for computing  $\pi$  can be deduced from the identity  $\pi/4 = 4 \arctan \frac{1}{5} - \arctan \frac{1}{239}$ . Determine the number of terms that must be summed to ensure an approximation to  $\pi$  to within  $10^{-3}$ .
6. Find the rates of convergence of the following sequences as  $n \rightarrow \infty$ .

a.  $\lim_{n \rightarrow \infty} \sin \frac{1}{n} = 0$

b.  $\lim_{n \rightarrow \infty} \sin \frac{1}{n^2} = 0$

c.  $\lim_{n \rightarrow \infty} \left( \sin \frac{1}{n} \right)^2 = 0$

d.  $\lim_{n \rightarrow \infty} [\ln(n+1) - \ln(n)] = 0$

7. Find the rates of convergence of the following functions as  $h \rightarrow 0$ .

a.  $\lim_{h \rightarrow 0} \frac{\sin h}{h} = 1$

b.  $\lim_{h \rightarrow 0} \frac{1 - \cos h}{h} = 0$

c.  $\lim_{h \rightarrow 0} \frac{\sin h - h \cos h}{h} = 0$

d.  $\lim_{h \rightarrow 0} \frac{1 - e^h}{h} = -1$

8. a. How many multiplications and additions are required to determine a sum of the form

$$\sum_{i=1}^n \sum_{j=1}^i a_i b_j?$$

- b. Modify the sum in part (a) to an equivalent form that reduces the number of computations.

9. Let  $P(x) = a_n x^n + a_{n-1} x^{n-1} + \cdots + a_1 x + a_0$  be a polynomial, and let  $x_0$  be given. Construct an algorithm to evaluate  $P(x_0)$  using nested multiplication.
10. Example 5 of Section 1.2 gives alternative formulas for the roots  $x_1$  and  $x_2$  of  $ax^2 + bx + c = 0$ . Construct an algorithm with input  $a, b, c$  and output  $x_1, x_2$  that computes the roots  $x_1$  and  $x_2$  (which may be equal or be complex conjugates) using the best formula for each root.
11. Construct an algorithm that has as input an integer  $n \geq 1$ , numbers  $x_0, x_1, \dots, x_n$ , and a number  $x$  and that produces as output the product  $(x - x_0)(x - x_1) \cdots (x - x_n)$ .
12. Assume that

$$\frac{1-2x}{1-x+x^2} + \frac{2x-4x^3}{1-x^2+x^4} + \frac{4x^3-8x^7}{1-x^4+x^8} + \cdots = \frac{1+2x}{1+x+x^2},$$

for  $x < 1$ , and let  $x = 0.25$ . Write and execute an algorithm that determines the number of terms needed on the left side of the equation so that the left side differs from the right side by less than  $10^{-6}$ .

13. a. Suppose that  $0 < q < p$  and that  $\alpha_n = \alpha + O(n^{-p})$ . Show that  $\alpha_n = \alpha + O(n^{-q})$ .  
 b. Make a table listing  $1/n, 1/n^2, 1/n^3$ , and  $1/n^4$  for  $n = 5, 10, 100$ , and  $1000$ , and discuss the varying rates of convergence of these sequences as  $n$  becomes large.
14. a. Suppose that  $0 < q < p$  and that  $F(h) = L + O(h^p)$ . Show that  $F(h) = L + O(h^q)$ .  
 b. Make a table listing  $h, h^2, h^3$ , and  $h^4$  for  $h = 0.5, 0.1, 0.01$ , and  $0.001$ , and discuss the varying rates of convergence of these powers of  $h$  as  $h$  approaches zero.
15. Suppose that as  $x$  approaches zero,

$$F_1(x) = L_1 + O(x^\alpha) \quad \text{and} \quad F_2(x) = L_2 + O(x^\beta).$$

Let  $c_1$  and  $c_2$  be nonzero constants, and define

$$F(x) = c_1 F_1(x) + c_2 F_2(x) \quad \text{and} \quad G(x) = F_1(c_1 x) + F_2(c_2 x).$$

Show that if  $\gamma = \text{minimum}\{\alpha, \beta\}$ , then as  $x$  approaches zero,

- a.  $F(x) = c_1 L_1 + c_2 L_2 + O(x^\gamma)$
  - b.  $G(x) = L_1 + L_2 + O(x^\gamma)$ .
16. The sequence  $\{F_n\}$  described by  $F_0 = 1$ ,  $F_1 = 1$ , and  $F_{n+2} = F_n + F_{n+1}$ , if  $n \geq 0$ , is called the *Fibonacci sequence*. Its terms occur naturally in many botanical species, particularly those with petals or scales arranged in the form of a logarithmic spiral. Consider the sequence  $\{x_n\}$ , where  $x_n = F_{n+1}/F_n$ . Assuming that  $\lim_{n \rightarrow \infty} x_n = x$  exists, show that  $x = (1 + \sqrt{5})/2$ . This number is called the *golden ratio*.
17. The Fibonacci sequence also satisfies the equation

$$F_n \equiv \tilde{F}_n = \frac{1}{\sqrt{5}} \left[ \left( \frac{1 + \sqrt{5}}{2} \right)^n - \left( \frac{1 - \sqrt{5}}{2} \right)^n \right].$$

- a. Write a Maple procedure to calculate  $F_{100}$ .
  - b. Use Maple with the default value of `Digits` followed by `evalf` to calculate  $\tilde{F}_{100}$ .
  - c. Why is the result from part (a) more accurate than the result from part (b)?
  - d. Why is the result from part (b) obtained more rapidly than the result from part (a)?
  - e. What results when you use the command `simplify` instead of `evalf` to compute  $\tilde{F}_{100}$ ?
18. The harmonic series  $1 + \frac{1}{2} + \frac{1}{3} + \frac{1}{4} + \dots$  diverges, but the sequence  $\gamma_n = 1 + \frac{1}{2} + \dots + \frac{1}{n} - \ln n$  converges, since  $\{\gamma_n\}$  is a bounded, nonincreasing sequence. The limit  $\gamma = 0.5772156649 \dots$  of the sequence  $\{\gamma_n\}$  is called Euler's constant.
- a. Use the default value of `Digits` in Maple to determine the value of  $n$  for  $\gamma_n$  to be within  $10^{-2}$  of  $\gamma$ .
  - b. Use the default value of `Digits` in Maple to determine the value of  $n$  for  $\gamma_n$  to be within  $10^{-3}$  of  $\gamma$ .
  - c. What happens if you use the default value of `Digits` in Maple to determine the value of  $n$  for  $\gamma_n$  to be within  $10^{-4}$  of  $\gamma$ ?

## 1.4 Numerical Software

Computer software packages for approximating the numerical solutions to problems are available in many forms. With this book, we have provided programs written in C, FORTRAN, Maple, Mathematica, MATLAB, Pascal, and Java that can be used to solve the problems given in the examples and exercises. These programs will give satisfactory results for most problems that you may need to solve, but they are what we call *special-purpose* programs. We use this term to distinguish these programs from those available in the standard mathematical subroutine libraries. The programs in these packages will be called *general purpose*.

The programs in general-purpose software packages differ in their intent from the algorithms and programs provided with this book. General-purpose software packages consider ways to reduce errors due to machine rounding, underflow, and overflow. They also describe the range of input that will lead to results of a certain specified accuracy.