

Chapter 2 Numerical Integration

2.1 Introduction

In some simple cases, the calculation of the definite integral

$$\int_a^b f(x)dx$$

(2.1.1)

is directly possible when the primitive (or antiderivative) function $F(x)$ is known

$$\int f(x)dx = F(x)$$

(2.1.2)

hence

$$\int_a^b f(x)dx = F(b) - F(a)$$

(2.1.3)

Most often, this is impossible and the only possible solution is numerical. Frequently, moreover, the function $f(x)$ is only known at a given number of points x_i , $i = 0, 1, \dots, n$. In this case, it is possible to search an approximation $g(x)$ of the function $f(x)$ and to proceed to a formal integration.

The interpolation polynomials $P_n(x)$ possess the required approximation properties and are easily integrable. Thus, they will be largely used in numerical integration (also called quadrature).

2.2 Newton and Cotes Closed Integration Formulas

The following integration formulas are called “closed” as they use the two basis points a and b to determine the approximation polynomial.

2.2.1 Global Integration on Interval $[a, b]$

Consider basis points uniformly distributed on interval $[a, b]$

$$x_i = a + ih, i = 0, 1, \dots, n \text{ with } h = \frac{b-a}{n}$$

(2.2.1)

Note that n is the degree of the interpolation polynomial $P_n(x)$ such that

$$P_n(x_i) = f(x_i) = f_i, i = 0, 1, \dots, n$$

(2.2.2)

For example, a Lagrange polynomial can be chosen as an interpolation polynomial. In this case

$$(2.2.3) \quad P_n(x) = \sum_{i=0}^n L_i(x) f_i$$

with

$$(2.2.4) \quad L_i(x) = \prod_{k=0, k \neq i}^n \frac{x - x_k}{x_i - x_k}$$

The variable $t \in [0, n]$ is introduced such that $x = a + ht$. The polynomial $L_i(x)$ becomes

$$(2.2.5) \quad L_i(x) = \phi_i(t) = \prod_{k=0, k \neq i}^n \frac{t - k}{i - k}$$

By integrating, we get

$$\int_a^b P_n(x) dx \left| \begin{array}{l} = \sum_{i=0}^n f_i \int_a^b L_i(x) dx \\ = h \sum_{i=0}^n f_i \int_0^n \phi_i(t) dt \\ = h \sum_{i=0}^n f_i w_i \end{array} \right.$$

(2.2.6)

The coefficients w_i are called weights; they depend only on n , thus they neither depend on the function f nor on the integration limits a and b . Recall that $h = (b - a)/n$.

Example: $n = 1$

$$w_0 = \int_0^1 \frac{t - 1}{0 - 1} dt = \int_0^1 (1 - t) dt = \frac{1}{2} w_1 = \int_0^1 \frac{t - 0}{1 - 0} dt = \int_0^1 t dt = \frac{1}{2}$$

(2.2.8)

which gives the following result:

$$(2.2.9) \quad \int_a^b P_1(x) dx = \frac{h}{2} (f_0 + f_1) = \frac{h}{2} [f(a) + f(b)]$$

corresponding to the trapezoidal rule with $h = (b - a)$ (Figure 2.2).
 Example: $n = 2$

$$w_0 = \int_0^2 \frac{t-1}{0-1} \frac{t-2}{0-2} dt = \frac{1}{2} \int_0^2 (t^2 - 3t + 2) dt = \frac{1}{3}$$

(2.2.10)

$$w_1 = \int_0^2 \frac{t-0}{1-0} \frac{t-2}{1-2} dt = - \int_0^2 (t^2 - 2t) dt = \frac{4}{3}$$

(2.2.11)

$$w_2 = \int_0^2 \frac{t-0}{2-0} \frac{t-1}{2-1} dt = \frac{1}{2} \int_0^2 (t^2 - t) dt = \frac{1}{3}$$

(2.2.12)

which gives the following result:

$$\int_a^b P_2(x) dx = \frac{h}{3} (f_0 + 4f_1 + f_2) = \frac{h}{3} [f(a) + 4f(\frac{a+b}{2}) + f(b)]$$

(2.2.13)

which is Simpson rule with $h = (b - a)/2$ (Figure 2.1).

By continuing, Table 2.1 results for different values of the degree n of the interpolation polynomial. From the degree n , the value of s results, then the weights w_i . The values σ_i are introduced only to display a table of integer values instead of fractional weights.

Newton-Cotes formulas thus give the approximation of the integral

$$\int_a^b P_n(x) dx = h \sum_{i=0}^n w_i f_i$$

(2.2.14)

with

$h = (b - a)/n$ (2.2.15) The weights w_i are such that their sum is equal to the degree n of interpolation polynomial

$\sum_{i=0}^n w_i = n$ (2.2.16) Let s be the lowest common denominator of the weights w_i . The integer numerators σ_i are such that

$$\sigma_i = s w_i \quad (2.2.17)$$

$$\sigma_i = \frac{n!}{(i+1)!(n-i)!} \quad (2.2.18)$$

$$\sigma_0 = 1$$

Fig. 2.2 Trapezoidal rule Example: For Simpson's rule, according to Table 2.1, we have $n = 2$, $n_s = 6$, thus $s = 3$. The weights result $w_0 = 1/3$, $w_1 = 4/3$, $w_2 = 1/3$. Newton-Cotes are now expressed as

$$\int_a^b P_n(x) dx = h \sum_{i=0}^n w_i f(x_i) = \frac{b-a}{n} \sum_{i=0}^n w_i f(x_i) \quad (2.2.18)$$

The error made by doing the numerical integration is equal to $\int_a^b P_n(x) dx - \int_a^b f(x) dx = h^{p+1} K f^{(p)}(\xi)$

(2.2.19) where $\xi \in [a, b]$. The values of the degree p and of the constant K only depend on the degree n of the interpolation polynomial. The error being of order p , any polynomial function of degree lower than p will be exactly integrated as the derivative of order p will be zero.

Numerical Methods and Optimization 49 2.2.2 Integration on Subintervals In general, Newton-Cotes formulas are not applied on all the interval $[a, b]$, but on the sequence of subintervals composing $[a, b]$. The type of subinterval depends on the order of the chosen method. The points x_i composing the interval $[a, b]$ are defined by

$$x_i = a + ih, \quad i = 0, 1, \dots, N \text{ with } h = \frac{b-a}{N} \quad (2.2.20)$$

It can be noticed that, in the previous formula, the definition of h is different from Equation (2.2.1). N must be chosen in agreement with the order n of the integration formula. • For the trapezoidal rule, a subinterval is defined by $[x_i, x_{i+1}]$. • For Simpson's rule, N is chosen even (the number of calculation points x_i is odd), a subinterval is defined by $[x_{2i}, x_{2i+1}, x_{2i+2}]$, $i = 0, 1, \dots, N/2 - 1$. • For the 3/8 rule, N is a multiple of 3 and a subinterval will be defined by $[x_{3i}, x_{3i+1}, x_{3i+2}, x_{3i+3}]$, $i = 0, 1, \dots, N/3 - 1$. • Application of the trapezoidal rule: On a subinterval, the trapezoidal rule gives $I_i = h/2 [f(x_i) + f(x_{i+1})]$ (2.2.21)

$$\begin{aligned} \text{Applying it to all the interval } [a, b], \text{ we get } I(h) &= \sum_{i=0}^{N-1} I_i = h \\ & \left[\frac{1}{2} f(a) + f(a+h) + \dots + f(b-h) + \frac{1}{2} f(b) \right] \\ &= \frac{b-a}{2N} \sum_{i=0}^{N-1} [f(a+ih) + f(a+(i+1)h)] \end{aligned}$$

(2.2.22) The function f is assumed to be continuously differentiable. On each subinterval, the error is equal to $I_i - \int_{x_i}^{x_{i+1}} f(x) dx = h^3/12 f''(\xi_i)$ (2.2.23)

Then the sum of the individual errors is $I(h) - \int_a^b f(x) dx = h^3/12 \sum_{i=0}^{N-1} f''(\xi_i)$ (2.2.24)

The summation term can be bounded

$$\sum_{i=0}^{N-1} f''(\xi_i) \leq N \max_{\xi \in [a, b]} f''(\xi) \quad (2.2.25)$$

As f'' is continuous, $[\min_{\xi \in [a, b]} f''(\xi), \max_{\xi \in [a, b]} f''(\xi)]$ such that

$$f^{(2)}(i) = \frac{1}{N-1} \sum_{i=0}^{N-1} f^{(2)}(i) \quad (2.2.26)$$

hence

$$I(h) = \int_a^b f(x) dx = \frac{b-a}{12} h^2 f^{(2)}(\xi), \quad [a, b] \quad (2.2.27)$$

This result means that the error done when a trapezoidal rule is used decreases like the square of h and thus the method is of order 2. • Application of Simpson's rule: With N even, on each subinterval $[x_{2i}, x_{2i+1}, x_{2i+2}]$, $i = 0, 1, \dots, N/2 - 1$. It gives $I_i = \frac{h}{3} [f(x_{2i}) + 4f(x_{2i+1}) + f(x_{2i+2})]$ with $h = \frac{b-a}{N}$

$$N \quad (2.2.28)$$

By summing these $N/2$ values, the approximation on $[a, b]$ results $I(h) = \frac{N}{2} \sum_{i=0}^{N/2-1} I_i = \frac{h}{3} [f(a) + 4f(a+h) + 2f(a+2h) + 4f(a+3h) + \dots + 2f(b-h) + 4f(b) + f(b)] = \frac{h}{3} [f(a) + f(b) + 2 \sum_{i=1}^{N/2-1} f(a+2ih) + 4 \sum_{i=0}^{N/2-1} f(a+(2i+1)h)]$

$$(2.2.29)$$

The error done is $S(h) = \int_a^b f(x) dx - \frac{h}{5} \sum_{i=0}^{N/2-1} f^{(4)}(i) = \frac{h^4}{90} \sum_{i=0}^{N/2-1} f^{(4)}(i) \quad (2.2.30)$ In the same way as for the trapezoidal rule, provided that f is 4 times continuously differentiable, it results that $S(h) = \int_a^b f(x) dx - \frac{b-a}{180} h^4 f^{(4)}(\xi) \quad (2.2.31)$

Thus Simpson's rule is a method of order 4.

2.3 Open Newton and Cotes Integration Formulas The following integration formulas are called "open" as they do not demand one or the other one of the bounds of the integration interval. The interpolation polynomial is of order $n-2$. Consider $n-1$ base points regularly spaced x_1, \dots, x_{n-1} . It is supposed that the lower integration limit a coincides with $x_0 = x_1 - h$ where h is the spacing between adjacent points. The upper limit b is not fixed. The integration formula is

Numerical Methods and Optimization 51

$$\int_a^b f(x) dx = \int_a^b P_{n-2}(x) dx \quad (2.3.1)$$

Thus, by defining

$$x = x_0 + h \quad (2.3.2)$$

we get for $n=2$

$$\int_{x_0}^{x_2} f(x) dx = 2h f(x_1) + \frac{h^3}{3} f^{(2)}(\xi) \quad (2.3.3)$$

$$= \frac{h^3}{3}$$

$$\int_{x_0}^{x_3} f(x) dx = \frac{3h^2}{2} [f(x_1) + f(x_2)] + \frac{3h^4}{4} f^{(2)}(\xi) \quad (2.3.4)$$

$$= \frac{4}{3} \int_{x_0}^{x_4} f(x) dx = \frac{4h^3}{3} [2f(x_1) + f(x_2) + 2f(x_3)] + \frac{14h^5}{45} f^{(4)}(\xi) \quad (2.3.5)$$

$$= \frac{5}{24} \int_{x_0}^{x_5} f(x) dx = \frac{5h^2}{24} [11f(x_1) + f(x_2) + f(x_3) + 11f(x_4)] + \frac{95h^5}{144} f^{(4)}(\xi) \quad (2.3.6)$$

The closed Newton-Cotes formulas are more accurate

than the open formulas as soon as a number of points larger than 2 or 3 is used. Thus, in general, it is better to use the closed formulas.

2.4 Conclusions on Newton and Cotes Integration Formulas The formulas with m points for m odd have the same order of accuracy as the formulas with $m + 1$ points. Their degree of precision is equal to m . A formula of degree of precision m exactly integrates all the polynomials of degree lower than or equal to m . The polynomials of larger degree are not exactly integrated. Thus, Simpson's rule exactly integrates polynomials of degree lower than or equal to 3. Except for the trapezoidal rule used because of its simplicity, it is preferable to use formulas with an odd number of base points than with an even number. The formulas with a number of base points larger than 8 are rarely used. Indeed, the rounding errors become large because of large weight factors with alternate signs. A way to reduce the error is the use of composite integration formulas. Rather than using a formula of a high order, it is often better to choose a formula having a low order, divide the integration interval $[a, b]$ into subintervals, and use the formula of low order separately on each subinterval.

52 Chapter 2. Numerical Integration 2.5 Repeated Integration by Dichotomy and Romberg's Integration Let $I_{N,1}$ be the estimation of the integral $\int_a^b f(x)dx$ (2.5.1) obtained by using the composite trapezoidal rule with a number n of subintervals such that $n = 2^N$. $I_{0,1}$ is the estimation of the integral obtained by using the simple trapezoidal rule (step = h)

$$I_{0,1} = (b - a) \frac{1}{2} [f(a) + f(b)]$$

(2.5.2) $I_{1,1}$ is the estimation of the integral obtained by using the simple trapezoidal rule applied two times (step = $h/2$)

$$I_{1,1} = (b - a) \frac{1}{2} [f(a) + f(b)] + f\left(\frac{a+b}{2}\right) \frac{(b-a)}{2}$$

$$= \frac{1}{2} [I_{0,1} + (b-a)f\left(\frac{a+b}{2}\right)]$$

(2.5.3) $I_{2,1}$ is the estimation of the integral obtained by using the simple trapezoidal rule applied four times (step = $h/4$)

$$I_{2,1} = (b - a) \frac{1}{2} [f(a) + f(b)] + 3 \sum_{i=1}^3 f\left(\frac{a+i(b-a)}{4}\right) \frac{(b-a)}{4}$$

The recurrence relation relating $I_{n,1}$ (step = $h/2^n$) to $I_{n-1,1}$ (step = $h/2^{n-1}$) is thus expressed as

$$I_{n,1} = \frac{1}{2^n} \left[I_{n-1,1} + (b-a) \sum_{i=1}^{2^{n-1}} f\left(a + i \frac{b-a}{2^n}\right) \right] \quad (2.5.5)$$

The error term corresponding to $I_{n,1}$ is equal to $(b-a)^3 \frac{1}{12} \frac{1}{2^{3n}} f''(\xi)$, $\xi \in [a, b]$ (2.5.6) Provided that the function $f(x)$ be continuous and bounded, $I_{n,1}$ converges to the exact value of the integral. Richardson's Extrapolation Now, introduce the general technique of Richardson's extrapolation. Given a quantity

Numerical Methods and Optimization 53 g_{approx} obtained by means of a discretization step h , g can be an integral, a derivative, ..., approximating the exact value g_{exact} . Suppose that the approximation is of order n , hence

$$g_{\text{exact}} = g_{\text{approx}}(h) + O(h^n) \quad (2.5.7)$$

which could be written as

$g_{\text{exact}} = g_{\text{approx}}(h) + a_n h^n + a_{n+1} h^{n+1} + O(h^{n+2})$ (2.5.8) where the coefficients a_i depend on the approximation method used. If instead of using a step h , we use $h/2$, Equation (2.5.8) becomes

$$g_{\text{exact}} = g_{\text{approx}}\left(\frac{h}{2}\right) + a_n \frac{h^n}{2^n} + a_{n+1} \frac{h^{n+1}}{2^{n+1}} + O(h^{n+2}) \quad (2.5.9)$$

where $g_{\text{approx}}\left(\frac{h}{2}\right)$

is in general a better approximation of g_{exact} than $g_{\text{approx}}(h)$. If Equations (2.5.8) and (2.5.9) are combined in order to eliminate the term about h^n , the approximation is then improved as the error will be at least about h^{n+1} . We thus get $(2^n - 1) g_{\text{exact}} = 2^n g_{\text{approx}}\left(\frac{h}{2}\right) - g_{\text{approx}}(h) + a_{n+1} \frac{2^n h^{n+1}}{2^n} + O(h^{n+2})$ (2.5.10)

Richardson's extrapolation formula results

$$g_{\text{exact}} = \frac{2^n g_{\text{approx}}\left(\frac{h}{2}\right) - g_{\text{approx}}(h)}{2^n - 1} + O(h^{n+1}) \quad (2.5.11)$$

demonstrating that this new formula gives a better result with an approximation of order $(n+1)$. Application of Richardson's Extrapolation Let us apply Richardson's extrapolation technique to a pair of adjacent elements of the sequence $I_{i,1}$ to obtain a better approximation of the integral. Each integral $I_{k,1}$ resulting from the trapezoidal rule is obtained with an error of order 2, thus $n = 2$. The application of Richardson's extrapolation then gives the relation

$$I = \frac{2^2 I_{k+1,1} - I_{k,1}}{2^2 - 1} + O(h^3) \quad (2.5.12)$$

as $I_{k+1,1}$ uses a step two times lower than $I_{k,1}$. The approximation results

$I = \frac{22}{3} I_{k+1,1} - I_{k,1}$ (2.5.13) Applied to the previous sequence, Richardson's extrapolation (Figure 2.3) gives for two adjacent elements

$$I_{k,2} = \frac{4I_{k+1,1} - I_{k,1}}{3} \quad (2.5.14)$$

54 Chapter 2. Numerical Integration For $k = 0$ corresponding to the trapezoidal rule applied once for $I_{0,1}$ on the interval $[a, b]$ and twice for $I_{1,1}$, Richardson's extrapolation yields

$$I_{0,2} = \frac{4I_{1,1} - I_{0,1}}{3} = \frac{(b-a)^6}{6} f(a) + 4f\left(\frac{a+b}{2}\right) + f(b)$$

that is Simpson's rule. The error on $I_{k,2}$ is equal to

$$\frac{(b-a)^5}{2880(2)^{4k}} f^{(4)}(\xi), \quad \xi \in [a, b] \quad (2.5.16)$$

$I_{1+1,1}$ (step = $h/2$) (step = $h/2 + 1$) Richardson 2 Fig. 2.3 Richardson's extrapolation By repeating Richardson's extrapolation, Romberg's extrapolation formula results

$$I_{k,j} = \frac{4^j I_{k+1,j} - I_{k,j}}{4^j - 1} \quad (2.5.17)$$

The error corresponding to $I_{k,j}$ is equal to $c(j) \frac{h^{2j}}{2^{2jk}} f^{(2j)}(\xi)$, $\xi \in [a, b]$ (2.5.18)

where $c(j)$ is a constant depending on a, b, j . It is interesting to note that if $I_{k,j}$ converges to the exact value of the integral when k increases, $I_{0,j}$ also converges when j increases. Thus, in the particular numerical Example 2.1, $I_{0,j}$ is close to the solution to about 10^6 for $j = 10$ whereas it is only the case for $I_{13,1}$. Example 2.1 : Integral by Richardson's extrapolation The following integral has been calculated: $I = \int_0^3 x \exp(x^2) dx$ (2.5.19)

by three different methods:

Numerical Methods and Optimization 55 (a) Gauss-Legendre quadrature with 5 points gives $I = 3963.45$ and with 10 points, it gives $I = 4051.04$. (b) Simpson's rule with 5 base points gives $I = 6441.21$, with 11 base points $I = 4258.18$, with 101 base points $I = 4051.07$. (c) Richardson's extrapolation gives the following Table 2.2. The values of the calculated integrals follow the notations of Equation (2.5.17) with tabulated values of $I_{k,j}$ obtained by iterative use of Richardson's extrapolation formula. It converges to 4051.042.

Table 2.2 Richardson's extrapolation. Values of $I_{k,j}$ are given $k \ j \ 1$
 2345678 1 36,463.878 12,183.089 6058.419 4295.426 4061.334 4051.171 4051.042
 4051.042 2 18,253.286 6441.211 4322.973 4062.248 4051.181 4051.042 4051.042
 3 9394.230 4455.363 4066.322 4051.224 4051.043 4051.042 4 5690.080 4090.637

4051.460 4051.043 4051.042 5 4490.498 4053.908 4051.050 4051.042 6 4163.056
4051.228 4051.042 7 4079.185 4051.054 8 4058.087

2.6 Numerical Integration with Irregularly Spaced Points Previously, all developed integration formulas were of the form

$\int_a^b f(x) dx \approx \sum_{i=0}^n w_i f(x_i)$ (2.6.1) where the $n + 1$ weights w_i are known from the $n + 1$ values x_i . When the points x_i are not fixed, there are $2n + 2$ unknowns (w_i and x_i), which allows us to determine a polynomial of degree $2n + 1$.

2.6.1 Reminder on Orthogonal Polynomials Two functions $g_m(x)$ and $g_n(x)$ belonging to a family of functions $g_i(x)$ are orthogonal with respect to a weight function $w(x)$ on the interval $[a, b]$ when, for all n

$$\int_a^b g_m(x) g_n(x) w(x) dx = 0 \text{ if } n \neq m \quad (2.6.2)$$

$$\int_a^b [g_n(x)]^2 w(x) dx = c(n) > 0 \quad (2.6.3)$$

56 Chapter 2. Numerical Integration where the notation $\int_a^b f(x) g(x) w(x) dx$ is called scalar product of the functions f and g relative to the weight function w . The scalar product is a number. Two functions are orthogonal when their scalar product is zero. A function is normalized when the scalar product of the function by itself is equal to 1. If all orthogonal functions two by two of the ensemble are normalized, the ensemble is orthonormal. In general, the value of c depends on n . A way to generate an ensemble of orthogonal polynomials for a given weight function $w(x)$ is to use the recurrence relation

$$P_{n+1}(x) = (x - a_n) P_n(x) - b_n P_{n-1}(x), \quad n = 0, 1, 2, \dots \quad (2.6.4)$$

with the coefficients defined by $a_n = \frac{\int_a^b x P_n(x) w(x) dx}{\int_a^b P_n(x)^2 w(x) dx}$, $n = 0, 1, 2, \dots$ $b_n = \frac{\int_a^b x P_{n-1}(x) P_n(x) w(x) dx}{\int_a^b P_{n-1}(x)^2 w(x) dx}$, $n = 1, 2, \dots$, any b_0

(2.6.5) To demonstrate Equation (2.6.5), it suffices to consider Equation (2.6.4) and to multiply by $w(x) P_n$ or $w(x) P_{n-1}$ respectively, then to take the integral of the new equation and to use the properties of orthogonal polynomials. The polynomials defined by (2.6.4) are monic, i.e. the coefficient of the monomial x^n of largest degree of $P_n(x)$ is equal to 1. If each polynomial is divided by $\sqrt{c(n)}$, the ensemble of polynomials becomes orthonormal. Other orthogonal polynomials can be met with different normalizations. Each polynomial $P_n(x)$ has exactly n distinct roots in the interval $[a, b]$. Among the known families of orthogonal functions, let us cite the family $(\sin kx)$ and the family $(\cos kx)$. The monomial functions are not orthogonal. On the opposite, there exist several families of orthogonal polynomials. Legendre polynomials: Legendre polynomials $P_n(x)$ are orthogonal on the interval $[-1, 1]$ with the unit weight function $w(x) = 1$.

$$P_m(x)P_n(x)dx = 0 \text{ if } n \neq m \quad (2.6.6)$$

and moreover

$$\int_{-1}^1 [P_n(x)]^2 dx = \frac{2}{2n+1} \quad (2.6.7)$$

The first Legendre polynomials are

Numerical Methods and Optimization 57

$$P_0(x) = 1 \quad P_1(x) = x \quad P_2(x) = \frac{1}{2}(3x^2 - 1) \quad P_3(x) = \frac{1}{2}(5x^3 - 3x) \quad P_4(x) = \frac{1}{8}(35x^4 - 30x^2 + 3) \dots \quad P_n(x) = \frac{2n-1}{n} x P_{n-1}(x) - \frac{n-1}{n} P_{n-2}(x) \quad (2.6.8)$$

Example 2.2 : Orthogonal Legendre polynomials To determine the orthogonal Legendre polynomials, consider the recurrence (2.6.4) $P_{n+1}(x) = (x - a_n)P_n(x) - b_n P_{n-1}(x)$, $n = 0, 1, 2, \dots$ (2.6.9)

with $P_0(x) = 1$, $P_1(x) = x$ and the equations of the coefficients (2.6.5), that is

$$\begin{aligned} a_n &= \frac{\int_{-1}^1 x P_n^2(x) dx}{\int_{-1}^1 P_n^2(x) dx}, \quad n = 1, 2, \dots \\ b_n &= \frac{\int_{-1}^1 x P_n(x) P_{n-1}(x) dx}{\int_{-1}^1 P_n^2(x) dx}, \quad n = 1, 2, \dots \end{aligned} \quad (2.6.10)$$

If $P_{n+1}(x)$ is calculated by Equation (2.6.9), we get a monic polynomial (whose coefficient of monomial x^{n+1} of highest degree is equal to 1) which does not satisfy Equation (2.6.7), thus we must set the Legendre polynomial equal to

$$P_{\text{Leg}n+1}(x) = c_{n+1} P_{n+1}(x) \quad (2.6.11)$$

where c_{n+1} is a coefficient such that

$$\int_{-1}^1 [P_{\text{Leg}n+1}(x)]^2 dx = \frac{2}{2(n+1)+1} = \int_{-1}^1 [c_{n+1} P_{n+1}(x)]^2 dx \quad (2.6.12)$$

hence

$$c_{n+1} = \sqrt{\frac{2}{2(n+1)+1}} \quad (2.6.13)$$

Thus, Table 2.3 results. Table 2.3 Orthogonal monic polynomials P and Legendre polynomials P_{Leg}

n	a_n	b_n	$P_{n+1}(x)$	$P_{\text{Leg}n+1}(x)$
0	0	0.3333	x^2	$0.3333 x^2$
1	0.5	0.2666	$x^3 - 0.6x$	$0.25x^3 - 0.375x$
2	0	0.2570	$x^4 - 0.8570x^2 + 0.0857$	$0.375x^4 - 0.75x^2 + 0.375$

58 Chapter 2. Numerical Integration Laguerre polynomials: Laguerre polynomials $L_n(x)$ are orthogonal on the interval $[0, +\infty[$ with the weight function $w(x) = \exp(-x)$

$$\int_0^{\infty} \exp(-x) L_m(x) L_n(x) dx = 0 \text{ if } n \neq m \quad (2.6.14)$$

$$\int_0^{\infty} \exp(-x) [L_n(x)]^2 dx = \frac{1}{n!} \quad (2.6.15)$$

The first Laguerre polynomials are $L_0(x) = 1$ $L_1(x) = x + 1$ $L_2(x) = x^2 - 2x + 2$ $L_3(x) = x^3 - 3x^2 + 6x - 6$... $L_n(x) = (2n - 1 - x)L_{n-1}(x) - (n - 1)^2 L_{n-2}(x)$

(2.6.16)

Chebyshev polynomials of the first kind: Chebyshev polynomials of the first kind $T_n(x)$ are orthogonal on the interval $[-1, 1]$ with the weight function $w(x) = 1/\sqrt{1 - x^2}$

$$\int_{-1}^1 \frac{1}{\sqrt{1 - x^2}} T_m(x) T_n(x) dx = 0 \text{ if } n \neq m$$

$$\int_{-1}^1 \frac{1}{\sqrt{1 - x^2}} [T_n(x)]^2 dx = \pi \text{ if } n = 0 \text{ if } n = 0$$

(2.6.17)

The first Chebyshev polynomials of the first kind are

$$T_0(x) = 1 \quad T_1(x) = x \quad T_2(x) = 2x^2 - 1 \quad T_3(x) = 4x^3 - 3x \quad \dots \quad T_n(x) = 2xT_{n-1}(x) - T_{n-2}(x)$$

(2.6.18)

Hermite polynomials: Hermite polynomials $H_n(x)$ are orthogonal on the interval $[-\infty, \infty]$ with the weight function $w(x) = \exp(-x^2)$ $\int_{-\infty}^{\infty} \exp(-x^2) H_m(x) H_n(x) dx = 0$ if $n \neq m$ (2.6.19) $\int_{-\infty}^{\infty} \exp(-x^2) [H_n(x)]^2 dx = \sqrt{\pi} 2^n n!$ (2.6.20)

The first Hermite polynomials are

Numerical Methods and Optimization 59

$$H_0(x) = 1 \quad H_1(x) = 2x \quad H_2(x) = 4x^2 - 2 \quad H_3(x) = 8x^3 - 12x \quad \dots \quad H_n(x) = 2xH_{n-1}(x) - (n - 1)H_{n-2}(x)$$

(2.6.21)

Each of these orthogonal polynomials of degree n with real coefficients has n distinct roots in its definition interval. Any polynomial of degree n can be represented as a linear combination of functions of any of the previous families.

2.6.2 Gauss–Legendre Quadrature

We estimate the integral in the same way as previously by integration of an approximation polynomial of degree n

$$\int_a^b f(x) dx \approx \int_a^b P_n(x) dx + \int_a^b R_n(x) dx \quad (2.6.22)$$

$R_n(x)$ being the error term. Let us use Lagrange interpolation polynomial $f(x) \approx \sum_{i=0}^n L_i(x) f(x_i) + R_n(x)$

$$\sum_{i=0}^n (x - x_i)$$

$$f(x) \approx \sum_{i=0}^n \frac{(x - x_i) \dots (x - x_{i-1})(x - x_{i+1}) \dots (x - x_n)}{(x_i - x_0) \dots (x_i - x_{i-1})(x_i - x_{i+1}) \dots (x_i - x_n)} f(x_i) \quad (2.6.23)$$

with

$$L_i(x) = \prod_{j=0, j \neq i}^n \frac{x - x_j}{x_i - x_j}$$

(2.6.24) The integration interval $[a, b]$ is transformed into $[-1, 1]$ by a change of variable

$$z = \frac{2x - (a + b)}{b - a} \quad (2.6.25)$$

and we define $F(z) = f(x)$, hence $F(z) = \sum_{i=0}^n Li(z)F(z_i) +$
 $\sum_{i=0}^n (z - z_i)$

$$F^{(n+1)}(\xi) \frac{(n+1)!}{(n+1)!}, \quad 1 \leq i \leq n \quad (2.6.26)$$

with

$$Li(z) = \sum_{j=0}^n j! \frac{z - z_j}{z_i - z_j} \quad (2.6.27)$$

Supposing that $f(x)$ is a polynomial of degree $2n+1$, then

60 Chapter 2. Numerical Integration

$F^{(n+1)}(\xi) \frac{(n+1)!}{(n+1)!} = q_n(\xi)$ (2.6.28) is a polynomial of degree n , belonging to the interval $[1, 1]$, but having no known value in this interval, to be able to pursue the calculations, $q_n(\xi)$ is transformed into a polynomial $q_n(z)$ on which it will be possible to work. An estimation of the integral is then given by

$$\int_{-1}^1 F(z) dz = \sum_{i=0}^n F(z_i) \int_{-1}^1 Li(z) dz = \sum_{i=0}^n w_i F(z_i) \quad (2.6.29)$$

with the weights $w_i = \int_{-1}^1 Li(z) dz = \int_{-1}^1 \sum_{j=0}^n j! \frac{z - z_j}{z_i - z_j} dz$

$\int_{-1}^1 dz$ (2.6.30) The integral to be calculated on $[a, b]$ is related to the integral calculated after change of variable by $\int_a^b f(x) dx = \int_a^b \frac{1}{2} \int_{-1}^1 F(z) dz = \int_a^b \sum_{i=0}^n w_i F(z_i) dz$ (2.6.31) Taking into account the previous remark about $q_n(\xi)$, the error term of the quadrature formula takes the form \int_{-1}^1

$$\sum_{i=0}^n (z - z_i) q_n(z) dz \quad (2.6.32)$$

The abscissas z_i must be chosen in order to minimize the error term. Consider the particular case of Gauss–Legendre quadrature. The polynomials $\sum_{i=0}^n (z - z_i)$ and $q_n(z)$ are expressed by means of Legendre polynomials P_i

$$\sum_{i=0}^n (z - z_i) = \sum_{i=0}^{n+1} b_i P_i(z) \quad (2.6.33)$$

$$q_n(z) = \sum_{i=0}^n c_i P_i(z) \quad (2.6.34)$$

The integral to minimize becomes

Numerical Methods and Optimization 61 \int_{-1}^1

$$\sum_{i=0}^n (z - z_i)$$

$$q_n(z) dz =$$

$$\begin{aligned} & \int_{-1}^1 \sum_{i=0}^n \sum_{j=0}^n b_i c_j P_i(z) P_j(z) + b_{n+1} \sum_{i=0}^n c_i P_i(z) P_{n+1}(z) dz \\ = & \int_{-1}^1 \sum_{i=0}^n b_i c_i [P_i(z)]^2 dz = \\ & \sum_{i=0}^n b_i c_i \int_{-1}^1 [P_i(z)]^2 dz \end{aligned} \quad (2.6.35)$$

as a result of orthogonality properties. The error term can be rendered equal to zero by imposing that the first $n + 1$ coefficients b_i be zero. There remains only one coefficient b_{n+1} different from zero so that, from Equation (2.6.33) $\sum_{i=0}^n (z - z_i) = b_{n+1} P_{n+1}(z)$ (2.6.36) As the coefficient of the term of degree n of the left-hand polynomial is equal to 1, it results that b_{n+1} is equal to the inverse of the coefficient of the term of highest degree of P_{n+1} . From the previous equality, it is obvious that the $n + 1$ base points z_i used in the integration formula are the $n+1$ roots of Legendre polynomial of degree $n+1$. Now that the base points are determined, the weights w_i can be calculated by Equation (2.6.30). Several different methods have been developed to efficiently calculate the weights. In the particular case of Legendre polynomials (Abramowitz and Stegun 1972), it is possible to use

$w_i = 2(1 - z_i^2)^{-1/2} [P'_{n+1}(z_i)]^{-2}$ (2.6.37) This constitutes Gauss–Legendre quadrature. Gauss–Legendre quadrature gives an exact integration result when the integrated function f is a polynomial of maximum degree $(2n + 1)$. The values of the roots z_i and the corresponding weights w_i for a family of given orthogonal polynomials are tabulated (Table 2.4). Remark: The calculation of the integral $I = \int_a^b f(x) dx$ (2.6.38) is brought back to the calculation of the integral approximated by the sum

$$I \approx \int_a^b f(z) dz \approx \sum_{i=0}^n w_i f(z_i) \quad (2.6.39)$$

62 Chapter 2. Numerical Integration Rather than making the change of variable $x \rightarrow z$ to determine the function $F(z)$ from $f(x)$, it is simpler to calculate the roots x_i on $[a, b]$ corresponding to the roots z_i on $[-1, 1]$ and to use the equality $f(x_i) = F(z_i)$. Thus, we get

$$I \approx \int_a^b f(x) dx \approx \sum_{i=0}^n w_i f(x_i) \quad (2.6.40)$$

Table 2.4 Gauss–Legendre quadrature formulas

Gauss–Legendre quadrature $\int_{-1}^1 F(z) dz \approx \sum_{i=0}^n w_i F(z_i)$

Roots z_i	Type	Weights w_i	$\pm 1/3$	Formula with 2 points ($n = 1$)	1.0000000000000000
0.0000000000000000	Formula with 3 points ($n = 2$)	$8/9 \pm 15/5$	$5/9 \pm 5/25$	525	
70 30/35	Formula with 4 points ($n = 3$)	$(18 + 30)/36$			
$\pm 5/25 + 70 30/35$	$(18 - 30)/36$	0.0000000000000000	Formula with 5 points ($n = 4$)	0.5688888888888888	± 0.538469310105683
0.478628670499366		± 0.906179845938664	0.236926885056189	± 0.238619186083197	Formula with 6 points ($n = 5$)
0.467913934572691		± 0.661209386466265	0.360761573048139	± 0.932469514203152	0.171324492379170
± 0.148874338981631	Formula with 10 points ($n = 9$)	0.295524224714753	± 0.433395394129247	0.269266719309996	± 0.679409568299024
0.219086362515982		± 0.865063366688985	0.149451349150581	± 0.973906528517172	0.066671344308688

Example 2.3 : Gauss–Legendre quadrature Calculate numerically the following integral:

Numerical Methods and Optimization 63

$I = \int_2^3 x^4 dx$ If it is integrated analytically, the exact value of this integral is $I=55$. We perform an integration according to Gauss–Legendre quadrature with 3 points ($n = 2$). We proceed in the following way: • Calculation of the abscissas x_i on interval $[2, 3]$ corresponding to the tabulated zeros z_i which are in $[-1, 1]$. • Calculation of the values of the function $f(x_i)$. • Evaluation of the integral according to Equation (2.6.40). Thus, Table 2.5 results. Table 2.5 Gauss–Legendre quadrature

z_i	x_i	$f(x_i)$	w_i	$w_i f(x_i)$
0	0.5	0.0625	8/9	0.05555
15/5	2.4364915	35.2419	5/9	19.5788
15/5	1.4364915	4.2580	5/9	2.36555

Finally

$I = \frac{5}{2} (0.05555 + 19.5788 + 2.36555) = 55$

Notice that, as the polynomial function to integrate had a degree lower than $(2n + 1)$ with $n = 2$, the integration is exact. Comparison with the trapezoidal rule and Simpson's rule without subintervals: Trapezoidal rule:

$I = \frac{5}{2} [1^2 (2) + 4 + 1^2 (3) + 4] = 242.5$

Simpson's rule:

$I = \frac{5}{2} [1^3 (2) + 4 + 4^3 (0.5) + 4 + 1^3 (3) + 4] = 81.04$

The interest of Gauss–Legendre quadrature is obvious with respect to the trapezoidal rule and Simpson's rule.

2.6.3 Gauss–Laguerre Quadrature

Gauss–Laguerre quadrature is based on the same principle as Gauss–Legendre quadrature. However, the integration formula takes into account the weight function under

the form $\int_0^\infty \exp(-z) F(z) dz = \sum_{i=0}^n w_i F(z_i) = \int_0^\infty G(z) dz = \sum_{i=0}^n w_i \exp(-z_i) G(z_i)$ (2.6.41)

64 Chapter 2. Numerical Integration where the z_i are the roots of Laguerre polynomial L_n and the weights are equal to

$w_i = \frac{(n!)^2}{(n+1)!} \frac{1}{[L_{n+1}(z_i)]^2}$ (2.6.42)

2.6.4 Gauss–Chebyshev Quadrature In the case of Chebyshev polynomials of the first kind, Gauss–Chebyshev quadrature gives $\int_{-1}^1 F(z) \sqrt{1-z^2} dz = \sum_{i=0}^n w_i F(z_i) = \int_{-1}^1 G(z) \sqrt{1-z^2} dz = \sum_{i=0}^n w_i G(z_i)$ (2.6.43) where the z_i are the roots of Chebyshev polynomial of the first kind T_n equal to

$z_i = \cos\left(\frac{(2i-1)\pi}{2n}\right)$
(2.6.44)

and the weights are equal to

$$w_i = \frac{1}{n} \quad (2.6.45)$$

2.6.5 Gauss–Hermite Quadrature Gauss–Hermite quadrature gives $\int_{-\infty}^{+\infty} \exp(-z^2) F(z) dz = \sum_{i=0}^n w_i F(z_i) = \int_{-\infty}^{+\infty} G(z) dz = \sum_{i=0}^n w_i \exp(-z_i^2) G(z_i)$ (2.6.46) where the z_i are the roots of Hermite polynomial H_n and the weights are equal to (by using the orthonormal ensemble of Hermite polynomials) $w_i = \frac{2}{[H'_n(z_i)]^2}$ (2.6.47)

2.7 Discussion and Conclusion Just as in approximation methods, the use of irregularly spaced points imposed by the quadrature method clearly demonstrates its advantage with respect to the accuracy of the result of numerical integration at the expense of a slightly more important work

Numerical Methods and Optimization 65 of understanding and design. Gauss–Legendre quadrature is the most frequently used. At a comparable level of precision, the number of calculations required to evaluate the integral is considerably lower than for the rules based on regularly spaced points. These latter are still used by many users who do not want to invest time, do a simple program, have the impression to master the method, in particular the trapezoidal rule. This shows the interest to use numerical libraries which are available and would provide the precision to them with a relatively reduced investment.

2.8 Exercise Set Exercise 2.8.1 (Easy) Calculate the following integral: $\int_0^2 x^4 dx$ (2.8.1) first by the trapezoidal method, then Simpson’s method, and finally Gauss–Legendre quadrature with 3 points ($n = 2$), by using in the three cases the bounds of the integration interval as calculation interval. Comment. Exercise 2.8.2 (Easy) The error function $\text{erf}(x)$ is defined by $\text{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x \exp(-t^2) dt = \frac{1}{\sqrt{\pi}} \int_{-\infty}^{+\infty} \exp(-t^2) dt$ (2.8.2) 1. Using Gauss–Legendre quadrature with 4 points, give an approximation of both integrals for $x = 1$. For the calculation of the second integral, it will be useful to think about the choice of the bound to use to replace $+$ and the influence of the value of the term $\exp(-t^2)$

2). It may be useful to make a few trials to better understand this influence. Discuss the results thus obtained. Deduce an approximation of $\text{erf}(x)$. 2. Do the same estimation with the first integral by Simpson’s method with a step $h = 0.1$. Compare the result thus obtained. Remark: The error function is used in many problems of physics. Consider a solid plate where the one-dimensional heat transfer is ruled by Fourier’s law

$T(x, t) = T_0 + \frac{q_0}{k} x^2$ (2.8.3) subjected to a constant temperature at the interface $T(x = 0, t) = T_s$. Let $T(x, t = 0) = T_0$ be the initial temperature. The

temperature along time (Incropera and DeWitt 1996) is expressed as

$$T(x, t) - T_s = T_0 - T_s = \text{erf} \left(\frac{x}{2\sqrt{\alpha t}} \right) \quad (2.8.4)$$

66 Chapter 2. Numerical Integration where α is the thermal diffusivity. Exercise 2.8.3 (Medium) The fugacity f (atm) of a gas at a pressure P (atm) and at temperature T is given by

$\ln f/P = \int_0^P \frac{Z-1}{P} dP$ (2.8.5) with the compressibility factor $Z = PV/(RT)$, R gas constant, and V molar volume (Smith et al. 2018; Vidal 1997). For methane, the experimental data of the compressibility factor Z with respect to pressure are given in Table 2.6. Table 2.6 Compressibility factor Z of methane with respect to pressure P

$P \text{ (atm)}$ 1 10 20 30 40 50 60 80 100 120 140 160 180 200
 Z 0.9940 0.9370 0.8683 0.7928 0.7034 0.5936 0.4515 0.3429 0.3767 0.4259 0.4753 0.5252 0.5752 0.6246

We desire to calculate the fugacity at $P = 200$ atm. 1. A first crude method would consist in using the trapezoidal method by using only the experimental data. Calculate in this way the fugacity with detailed calculation. 2. A second method would consist in using Gauss–Legendre quadrature. For that purpose, we propose an approximation function of the form $Z = 1 + 0.00858 P - 0.000463 P \ln(P + 1) + 0.0000475 P^2$ (2.8.6) By explaining the steps of the calculation, without fully explaining Gauss–Legendre quadrature, calculate the fugacity. Exercise 2.8.4 (Easy) Calculate the following integral: $I = \int_0^1 \frac{1}{1+x^4} dx$ (2.8.7)

Numerical Methods and Optimization 67 by both methods, 1. Simpson’s method with the step $h = 0.25$. 2. Gauss–Legendre quadrature with five points. Remark: To calculate this integral, it is recommended to divide the integration domain as $[0, 1]$ and $[1, +\infty[$, and then to do a change of variable $t = 1/x$ on the second domain. Exercise 2.8.5 (Medium) Calculate the following integral: $I = \int_0^1 \int_0^1 \frac{x^2 + y^2}{1+x^2+y^2} dx dy$ (2.8.8) by Gauss–Legendre quadrature with three points, clearly explaining the technique used and giving intermediate results. Exercise 2.8.6 (Easy) Calculate the following integral: $I = \int_0^2 2 \exp(x^2) dx$ (2.8.9)

by 5-point Gauss–Legendre quadrature.

References M. Abramowitz and I. A. Stegun. Handbook of Mathematical Functions with Formulas, Graphs and Mathematical Tables. Dover, New York, 1972. F. P. Incropera and D. P. DeWitt. Fundamentals of Heat and Mass Transfer. John Wiley, New York, 4th edition, 1996. J. M. Smith, H. C. Van Ness, M. M. Abbott, and M. T. Swihart. Introduction to Chemical Engineering Thermodynamics. McGraw-Hill, New York, 8th edition, 2018.

J. Vidal. Thermodynamique. Technip, Paris, 1997.

Chapter 3 Equation Solving by Iterative Methods

3.1 Introduction The problem is to develop adequate methods to find the solutions of the general equation

$f(x) = 0$ (3.1.1) The roots will be noted i . In a given number of cases, $f(x)$ will be supposed to be a polynomial of degree n

$$f(x) = x^n + a_1 x^{n-1} + \dots + a_{n-1} x + a_n \quad (3.1.2)$$

but some methods are applicable to any type of function. There exist four main classes of methods (Gritton et al. 2001) to find the roots of a nonlinear equation of the form

$$f(x) = 0 \quad (3.1.3)$$

1. Local methods that require an initial estimation of the root (e.g. successive substitutions, Newton). Frequently, local methods are designed to search for only one

real root of the nonlinear equation, even if multiple roots exist. Nevertheless, they are often very robust and nearly always converge (e.g. Newton, quasi-Newton), but they present the drawback to need to provide an initial estimation sufficiently close to the root. 2. Global methods that find a root from an arbitrary initial value (e.g. homotopy). They are adapted to the search of multiple roots. 3. Interval methods that find all the roots in a specified domain of x (e.g. dichotomy, regula falsi). They are robust but slow. 4. Graphical methods or spreadsheet that uses a graphical view of $f(x)$ in a specified domain of x . Some of the methods presented below (Graeffe, Bernoulli, Bairstow) are more interesting from a mathematical point of view than for real applications, but they present a historical interest and their exposure may promote future ideas. © The Author(s), under exclusive license to Springer Nature Switzerland AG 2021 J.-P. Corriou, Numerical Methods and Optimization, Springer Optimization and Its Applications 187, https://doi.org/10.1007/978-3-030-89366-8_3

69

70 Chapter 3. Equation Solving by Iterative Methods 3.2 Graeffe's Method Graeffe's method is a global method, as it gives a simultaneous approximation of all roots. Consider a monic polynomial f of type (3.1.2). To f , the following adjoint function is associated

$$f_1(x) = (1 - x_1^n) f(x) \quad f_2(x) = (x_1^2 - 2x_1 + 1)(x_2^2 - 2x_2 + 1) \dots (x_n^2 - 2x_n + 1) \quad (3.2.1)$$

where x_i are the searched roots, ordered by decreasing modulus.

As $f_1(x)$ contains only even powers, a new function is defined $f_2(x) = (f_1(x) - x_1^{2n}) / (x_1^2 - 2x_1 + 1) \dots (x_n^2 - 2x_n + 1)$ (3.2.2) which has the property that its roots are

the squares of the roots of f . The operation can be repeated, and we obtain a sequence of polynomials $f_2, f_4, f_8 \dots$ such that

$$f_m(x) = (x - m_1)(x - m_2) \dots (x - m_n) \quad (3.2.3)$$

where m is an integer positive of 2 and f_m has roots m_1, m_2, \dots, m_n . The aim of this sequence is to form an equation whose roots have very different orders of magnitude, that is, if the roots are real, the ratios m_i/m_j can be made as small as desired when

$$\begin{aligned} f_m(x) &= x^n (m_1 + \dots) x^{n-1} + (m_1 m_2 + \dots) x^{n-2} \\ &+ (m_1 m_2 m_3 + \dots) x^{n-3} + \dots + (1) x^n (m_1 m_2 \dots m_n) \\ &= x^n A_1 x^{n-1} + \dots + (1) x^n A_n \\ &= x^n A_1 x^{n-1} + \dots + (1) x^n A_n \end{aligned} \quad (3.2.4)$$

the approximations result $m_1 = A_1, m_2 = A_2/A_1, \dots, m_n = A_n/A_{n-1}$ (3.2.5) hence an approximation of the absolute values or of the moduli of the searched roots by taking the m th root. The sign of the roots is not determined by this method. It must be verified by substitution in the original equation. If multiple or complex roots exist, $m_i = m_{i+1}$, the equation

$$A_i x^2 - A_i x + A_{i+1} = 0 \quad (3.2.6)$$

gives approximations of m_i and m_{i+1} .

Graeffe's method presents some numerical difficulties and is not commonly used. Example 3.1 : Graeffe's method Consider a polynomial with real roots, equal to 1, 2, 3, and 4. This polynomial is equal to

$$P(x) = x^4 + 2x^3 - 13x^2 - 14x + 24 \quad (3.2.7)$$

Table 3.1 Graeffe's method: root finding for a polynomial having real roots

m_1	A_1/m_1	(A_2/A_1)	$1/m_1$	(A_3/A_2)	$1/m_1$	(A_4/A_3)	$1/m_1$	$1/2$	5.4772	3.0166	1.7331	0.8381	$2/4$	4.3376	2.9409
1.9173	0.9812	$3/8$	4.0497	2.9787	1.9904	0.9994	$4/16$	4.0024	2.9984	1.9998	0.9999	$5/32$	4.0000	3.0000	2.0000
1.0000															

Using Graeffe's method, the successive polynomials are found $f_2(x) = x^4 - 30x^3 + 273x^2 - 820x + 576$ $f_4(x) = x^4 - 354x^3 + 26481x^2 - 357904x + 331776$ $f_8(x) = x^4 - 72354x^3 + 448510881x^2 - 110523752704x + 110075314176$

(3.2.8) This shows that the polynomial coefficients increase very rapidly. In Table 3.1, the values of A_1/m_1 ,

$(A_2/A_1)^{1/m}, \dots$, have been gathered to highlight the limits that are the ordered root moduli. The convergence is fast. However, it must be noted that the coefficients A_i very rapidly take very large values, which poses huge numerical problems. In the case of a polynomial with complex roots, Graeffe's method is difficult to use. At the best, it allows to have an estimation of m_i .

3.3 Bernoulli's Method Bernoulli's method to find a root k of the polynomial

$P(x) = \sum_{i=0}^n a_i x^i$ with $a_0 = 1$ (3.3.1) first consists of building a sequence u_i by associating to each monomial x^k a term u_i^k (thus $i \leq n$). To understand the interest of the building of the sequence u_i , first consider the fact that the roots i of the polynomial $P(x)$ are supposed to be ordered according to their modulus $|i| \leq \dots \leq |i_n|$. Express that i is a root $a_0 i^n + a_1 i^{n-1} + \dots + a_{n-1} i + a_n = 0$

(3.3.2)

72 Chapter 3. Equation Solving by Iterative Methods Each of the previous equations is then multiplied by an arbitrary coefficient c_i , we sum the rows, that is $c_1(a_0 i^{n-1} + a_1 i^{n-2} + \dots + a_{n-1} i + a_n) + \dots + c_n(a_0 i^n + a_1 i^{n-1} + \dots + a_{n-1} i + a_n) = 0$ (3.3.3)

and then we order again with respect to the coefficients a_i , i.e. $a_0(c_1 i^{n-1} + \dots + c_n i^n) + a_1(c_1 i^{n-2} + \dots + c_n i^{n-1}) + \dots + a_n(c_1 + \dots + c_n) = 0$ (3.3.4)

We set

$$u_i = c_1 i^{n-1} + \dots + c_n i^n$$

$$i = 0, 1, \dots, n \quad (3.3.5)$$

hence

$$a_0 u_n + a_1 u_{n-1} + \dots + a_n u_0 = 0 \quad (3.3.6)$$

To simplify the writing, we consider $a_0 = 1$, hence $u_n = a_1 u_{n-1} + \dots + a_n u_0$
 $= \sum_{i=1}^n a_i u_{i-1}$ (3.3.7)

By extension, the sequence is defined $u_i = \sum_{j=1}^n a_j u_{i-j}$ for $i \leq n$ (3.3.8)
 Thus, it can be seen that to define this sequence, it suffices to arbitrarily choose the real numbers u_i for $0 \leq i \leq n-1$. From the form of the solutions u_i previously given

$$u_i = c_1 i^{n-1} + \dots + c_n i^n$$

$$i = 0, 1, \dots, n-1 \quad (3.3.9)$$

a system of n equations with n unknowns c_i

j results (Vandermonde determinant). To

find the solution, an initial vector u_i as simple as possible can be chosen

$$u_i = 0, \quad 0 \leq i \leq n-1 \quad \text{and} \quad u_n = 1 \quad (3.3.10)$$

Equation (3.3.8) can also be written as

$$u_{n+k} = a_1 u_{n+k-1} \cdots a_n u_k \quad (3.3.11)$$

Now consider the ratio of two successive terms

$$\frac{u_{n+k}}{u_{n+k-1}} = a_1 \frac{u_{n+k-1}}{u_{n+k-2}} \cdots a_n \frac{u_k}{u_{k-1}} \quad (3.3.12)$$

and suppose that this ratio admits a limit l when $k \rightarrow \infty$

$$\lim_{k \rightarrow \infty} \frac{u_{n+k}}{u_{n+k-1}} = l \quad (3.3.13)$$

From Equation (3.3.12), we deduce

$$l^{n+1} = a_1 l^n \cdots a_n l \quad (3.3.14)$$

hence the limit l is a root of the polynomial. Theorem:

If l_1, l_2, \dots, l_j are the roots of the polynomial, $l = \lim_{i \rightarrow \infty} \frac{u_i}{u_{i-1}}$, (3.3.15) thus the limit l corresponds to the root of larger modulus. Bernoulli's method remains valid even when multiple roots exist: $l_1 = l_2 = \dots = l_j$ provided we have: $l_j \neq l_{j+1}$. If there are no multiple roots, the unique solution for $i \geq 0$ is of the form

$$u_i = c_1 l_1^i + \dots + c_n l_n^i$$

$n, 0 \leq i \leq n-1$ (3.3.16) the values of the coefficients c_i depending on the initial value taken for u_i . In the case where a double root exists: $l_1 = l_2$ with $l_1 \neq l_3$, the solution is of the form

$$u_i = c_1 l_1^i + i c_2 l_1^i + c_3 l_3^i + \dots + c_n l_n^i$$

n (3.3.17) To find the solution, among all possible sequences, it may be useful to associate the two following sequences v_i and t_i to the sequence u_i

$$v_i = u_{i+1} - u_i \quad (3.3.18)$$

and

$$t_i = u_i u_{i+1} - u_{i+1} u_{i+2} \quad (3.3.19)$$

Thus, in the case (in particular, if l_1 and l_2 are complex) where

$$l_1 = l_2 \neq l_3 \quad (3.3.20)$$

we obtain the following results $\lim_{i \rightarrow \infty} \frac{v_i}{v_{i-1}} = l_1$ (3.3.21)

and

$$\lim_{i \rightarrow \infty} \frac{t_{i+1}}{t_i} = l_1^2 \quad (3.3.22)$$

Bernoulli's method is rarely employed since the intensive use of computers. Example 3.2 : Bernoulli's method Consider a polynomial with real roots, equal to 1, 2, 3, and 4. This polynomial is $P(x) = x^4 - 2x^3 - 13x^2 - 14x + 24$ (3.3.23) First, the sequence u_i is calculated by initializing with $u_0 = 0, u_1 = 0, u_2 = 0, u_3 = 1$. The other terms of this sequence obey Equation (3.3.16). The sequences v_i and t_i are also calculated according to Equations (3.3.18) and (3.3.19), respectively. The value of 1 results as the root of largest

1.4228 2.3650 2.5732 100 1.8325 2.3912 2.5670

3.4 Bairstow's Method Bairstow's method allows us to find the complex conjugate roots of a real polynomial by noticing that they correspond to the factorization of a real polynomial of degree 2. To calculate these roots by Newton's method, it would be necessary to use numerical complex calculation. In a general manner, a polynomial $P(x) = \sum_{i=0}^n a_i x^i$ (3.4.1)

by setting $P_0(x) = P(x)$ can be written under the form

$P_0(x) = P_1(x)(x^2 - rx - q) + A_0x + B_0$ (3.4.2) where the polynomial $P_1(x)$ is of degree $n-2$ and the remainder of the division of $P_0(x)$ by $P_1(x)$ is equal to $A_0x + B_0$. The coefficients A_0 and B_0 depend on the values of r and q ; thus the issue is to find the values of these coefficients that make $A_0(r, q)$ and $B_0(r, q)$ equal to zero. Indeed, Bairstow's method makes use of Newton-Raphson method that will be examined in Section 5.12, in the solution of systems of nonlinear equations. The recurrence relation giving r and q is

76 Chapter 3. Equation Solving by Iterative Methods

$r_{i+1} \quad q_{i+1}$

=

$r_i \quad q_i$

$A_0(r_i, q_i) \quad B_0(r_i, q_i)$

$r_{i+1} = r_i \quad q_{i+1} = q_i$

$A_0(r_i, q_i) \quad B_0(r_i, q_i)$

(3.4.3)

However, $A_0(r, q)$ and $B_0(r, q)$ are unknown, thus the elements of the Jacobian matrix also. Consequently, it is necessary to determine those four partial derivatives that are present in the following identities: $P_0(x) - r = (x^2 - rx - q)P_1(x) + A_0x + B_0 = 0$ (3.4.4)

$P_0(x) - q = (x^2 - rx - q)P_1(x) + A_0x + B_0 = 0$ (3.4.5) After this first division by $(x^2 - rx - q)$, the operation can be repeated; hence $P_1(x) = P_2(x)(x^2 - rx - q) + A_1x + B_1$ (3.4.6)

Supposing that $(x^2 - rx - q = 0)$ has two distinct roots x_0, x_1 , we get

$P_1(x_i) = A_1x_i + B_1, i = 0, 1$ (3.4.7)

and both identities (3.4.4) and (3.4.5) become

$x_i(A_1x_i + B_1) + A_0r x_i + B_0r = 0$

$(A_1x_i + B_1) + A_0q x_i + B_0q = 0, i = 0, 1$ (3.4.8) From the second equation of Equation (3.4.8) and from Equation (3.4.7), we draw

$A_0q = A_1, B_0$

$q = B_1$ (3.4.9)

and consequently the first equation of (3.4.8) becomes

$$x^2 \sum_{i=0}^n A_i q + \sum_{i=0}^n (A_i r - B_i q) + B_0 r = 0, \quad i = 0, 1 \quad (3.4.10)$$

$$\text{As } x^2 \sum_{i=0}^n A_i = r \sum_{i=0}^n A_i + q, \text{ it gives } \sum_{i=0}^n (A_i r - B_i q - A_i q r) + B_0 r - A_0 q q = 0, \quad i = 0, 1 \quad (3.4.11)$$

hence

$$A_0 r - B_0 q - A_0 q r = 0 \quad (3.4.12) \quad B_0 r - A_0 q q = 0 \quad (3.4.13)$$

Numerical Methods and Optimization 77 and finally

$$A_0 q = A_1, \quad B_0 q = B_1$$

$$A_0 r = r A_1 + B_1, \quad B_0 r = q A_1$$

$$(3.4.14)$$

By taking

$$P_0(x) = \sum_{i=0}^n a_i x^i \text{ and } P_1(x) = \sum_{i=0}^n b_i x^{2i} \quad (3.4.15)$$

$$\text{it results } \sum_{i=0}^n a_i x^i = (x^2 - r x - q) \sum_{i=0}^n b_i x^{2i} + A_0 x + B_0 \quad (3.4.16)$$

By identifying with respect to the successive powers of x by decreasing order, we get the relations $a_0 = b_0$ $a_1 = b_1 - r b_0$ $a_2 = b_2 - r b_1 - q b_0 \dots a_i = b_i - r b_{i-1} - q b_{i-2}$ $\dots a_{n-1} = r b_{n-2} - q b_{n-3} + A_0$ $a_n = q b_{n-2} + B_0$

$$(3.4.17)$$

and the values of A_0 and B_0 result by means of Horner's scheme

$$b_0 = a_0 \quad b_1 = b_0 r + a_1 \quad b_2 = b_0 q + b_1 r + a_2 \dots b_i = b_{i-2} q + b_{i-1} r + a_i, \quad i = 2, \dots, n-2 \dots A_0 = b_{n-3} q + b_{n-2} r + a_{n-1} \quad B_0 = b_{n-2} q + a_n$$

$$(3.4.18)$$

The process is repeated with P_1 and so on until exhaustion. Example 3.4 : Bairstow's method Consider the following real polynomial

$$P(x) = x^4 + 2x^3 + 3x^2 + 4x + 5 \quad (3.4.19)$$

having complex conjugate roots, equal to $0.2878 \pm 1.4160i$ and $1.2878 \pm 0.8578i$. By applying Bairstow's method, the values of $A_0(r, q)$ and $B_0(r, q)$ are found by applying Equation (3.4.18)

$$\begin{aligned} 78 \text{ Chapter 3. Equation Solving by Iterative Methods } & b_0 = 1 \quad b_1 = r + 2 \\ b_2 = q + r(r + 2) + 3 & A_0 = (r + 2)q + (q + (r + 2)r + 3)r + 4 = 2qr \\ + 2q + r^3 + 2r^2 + 3r + 4 & B_0 = (q + (r + 2)r + 3)q + 5 = q^2 + qr^2 + 2qr + 3q + 5 \end{aligned}$$

$$(3.4.20)$$

Applying Equation (3.4.3), Newton-Raphson's algorithm follows at iteration i

$$r, q$$

$$i+1 =$$

$$r, q$$

$$i$$

$$2q + 3r^2 + 4r + 32r + 22qr + 2q^2q + r^2 + 2r + 31i$$

$$A_0 B_0$$

i (3.4.21) The initialization is simply done with $(r, q) = (0, 0)$. After a limited number of iterations, Table 3.4 results. Table 3.4 Bairstow's method: roots solving for a real polynomial having complex conjugate roots Iteration

ir	q	A0	B0	1	00	4	5	2	0.2222	1.6667	0.8285	3.4362	3	1.0132	1.3463	4.7119
1.3364	4	0.5601	1.6794	1.2433	0.374	5	0.5514	2.0700	0.0071	0.1630	6	0.5760	2.0880	0.0013	0.0023	7

0.5756 2.0882 0.62 $\times 10^6$ 0.029 $\times 10^6$

Having found $(A_0, B_0)(0, 0)$, the searched polynomial results

$x^2 - r x - q = x^2 - 0.5756 x + 2.0882$ (3.4.22) whose roots are $0.2878 \pm 1.4160i$. In a general way, then it would be necessary to continue the method by applying it to the polynomial P_1 obtained by exact division of $P(x)$ by $(x^2 - r x - q)$. In the present case, as $P_1(x)$ is of degree 2, the solution is immediate.

3.5 Existence of a Root of a Function Before searching for the root of a continuous function f by an iterative method such that $f() = 0$, it is essential to find an interval $[a, b]$ such that

$f(a) f(b) < 0$ (3.5.1) This guarantees the existence of a solution in the interval $[a, b]$. Finding the interval $[a, b]$ is the initial step of bracketing.

To compare different root-finding methods, the following equation will be considered

$$f(x) = \exp(x) - x^2$$