



Numerical Differentiation of Analytic Functions

Author(s): J. N. Lyness and C. B. Moler

Source: *SIAM Journal on Numerical Analysis*, Vol. 4, No. 2 (Jun., 1967), pp. 202-210

Published by: Society for Industrial and Applied Mathematics

Stable URL: <http://www.jstor.org/stable/2949389>

Accessed: 14/12/2009 10:08

Your use of the JSTOR archive indicates your acceptance of JSTOR's Terms and Conditions of Use, available at <http://www.jstor.org/page/info/about/policies/terms.jsp>. JSTOR's Terms and Conditions of Use provides, in part, that unless you have obtained prior permission, you may not download an entire issue of a journal or multiple copies of articles, and you may use content in the JSTOR archive only for your personal, non-commercial use.

Please contact the publisher regarding any further use of this work. Publisher contact information may be obtained at <http://www.jstor.org/action/showPublisher?publisherCode=siam>.

Each copy of any part of a JSTOR transmission must contain the same copyright notice that appears on the screen or printed page of such transmission.

JSTOR is a not-for-profit service that helps scholars, researchers, and students discover, use, and build upon a wide range of content in a trusted digital archive. We use information technology and tools to increase productivity and facilitate new forms of scholarship. For more information about JSTOR, please contact support@jstor.org.



Society for Industrial and Applied Mathematics is collaborating with JSTOR to digitize, preserve and extend access to *SIAM Journal on Numerical Analysis*.

<http://www.jstor.org>

NUMERICAL DIFFERENTIATION OF ANALYTIC FUNCTIONS*

J. N. LYNESS† AND C. B. MOLER‡

1. Introduction. Many algebraic computer languages now include facilities for the evaluation of functions of a complex variable. Such facilities can be effectively used for numerical differentiation. The method we describe is appropriate for computing the derivatives $f^{(n)}(x)$ of any analytic function which can be evaluated at points in the complex plane near x . Such a function might be a complicated rational combination of elementary functions of x which is difficult to differentiate analytically.

In this paper we derive several formulas, any of which is suitable for evaluating the n th derivative of a complex analytic function at a point in terms of function evaluations at neighboring points. For simplicity we assume throughout this paper that the derivatives are to be evaluated at the origin. In this Introduction we state one formula; in succeeding sections we prove this and related formulas, discuss their degree and the error in their application, and give an example.

A trapezoidal rule sum operator $R^{[m,1]}$ may be defined as follows:

$$(1.1) \quad R^{[m,1]}g(t) = \frac{1}{m} \left[\frac{1}{2}g(0) + g\left(\frac{1}{m}\right) + \cdots + g\left(\frac{m-1}{m}\right) + \frac{1}{2}g(1) \right].$$

In (2.12) below we define the Möbius numbers μ_n , $n = 1, 2, \dots$, which have the values $+1, 0$, or -1 . In terms of

$$(1.2) \quad g(r; t) = g(t) = \operatorname{Re} f(re^{2\pi it}),$$

one of our results may be stated as follows.

THEOREM 1. *If $f(x)$ is a real function of x and $f(z)$ is analytic in a neighborhood that contains the circle $|z| = r$, then*

$$(1.3) \quad a_n = \frac{f^{(n)}(0)}{n!} = \frac{1}{r^n} \sum_{m=1}^{\infty} \mu_m [R^{[mn,1]}g(t) - f(0)].$$

In other words, the n th derivative of f can be computed from the trapezoidal rule operating on the real part of f . The only complex arithmetic required is in the computation of this real part.

This result is based on Cauchy's theorem, which expresses the n th deriva-

* Received by the editors August 15, 1966, and in revised form December 20, 1966. This research was sponsored in part by the United States Atomic Energy Commission under contract with the Union Carbide Corporation.

† Mathematics Division, Oak Ridge National Laboratory, Oak Ridge, Tennessee 37830. On leave from the University of New South Wales, Sidney, Australia.

‡ Department of Mathematics, University of Michigan, Ann Arbor, Michigan.

tive of an analytic function in terms of a contour integral. This integral contains an oscillating component and may be evaluated using the Poisson summation formula and the Möbius inversion of series.

Since this method is based on numerical quadrature, it does not show the sensitivity to roundoff error in the function evaluations that is characteristic of finite difference methods.

2. Proof of Theorem 1. Cauchy's theorem may be stated as

$$(2.1) \quad a_n \equiv \frac{f^{(n)}(0)}{n!} = \frac{1}{2\pi i} \int_C \frac{f(z)}{z^{n+1}} dz, \quad n = 0, 1, 2, \dots,$$

where $f(z)$ is analytic in a neighborhood which includes a closed contour C which surrounds the origin.

If we choose C to be the circle $|z| = r$ and use variable t defined by

$$(2.2) \quad z = re^{2\pi it},$$

we find by direct substitution

$$(2.3) \quad a_0 = f(0) = \int_0^1 f(re^{2\pi it}) dt,$$

$$(2.4) \quad a_n = \frac{1}{r^n} \int_0^1 f(re^{2\pi it}) e^{-2\pi int} dt, \quad n = 1, 2, 3, \dots$$

Using (2.3) with $f(z)$ replaced by $z^n f(z)$, we find

$$(2.5) \quad 0 = \frac{1}{r^n} \int_0^1 f(re^{2\pi it}) e^{2\pi int} dt, \quad n = 1, 2, 3, \dots$$

Equations (2.4) and (2.5) may be combined to give either of the two forms:

$$(2.6) \quad a_n = \frac{2}{r^n} \int_0^1 f(re^{2\pi it}) \cos 2\pi nt dt,$$

$$(2.7) \quad a_n = \frac{-2i}{r^n} \int_0^1 f(re^{2\pi it}) \sin 2\pi nt dt.$$

The Poisson summation formula in its finite form [5] connects the trapezoidal rule sum operator (1.1) with the Fourier coefficients (2.6). It states

$$(2.8) \quad \begin{aligned} R^{[n,1]} f(re^{2\pi it}) &= \sum_{k=-\infty}^{\infty} \int_0^1 f(re^{2\pi it}) e^{2\pi iknt} dt \\ &= \int_0^1 f(re^{2\pi it}) dt + 2 \sum_{k=1}^{\infty} \int_0^1 f(re^{2\pi it}) \cos 2\pi knt dt. \end{aligned}$$

If we denote by b_n the error in the n -point trapezoidal rule:

$$(2.9) \quad \begin{aligned} b_n &= \left[R^{[n,1]}f(re^{2\pi i t}) - \int_0^1 f(re^{2\pi i t}) dt \right] \\ &= R^{[n,1]}f(re^{2\pi i t}) - f(0), \end{aligned}$$

then (2.8) may be written in the form

$$(2.10) \quad b_n = r^n a_n + r^{2n} a_{2n} + r^{3n} a_{3n} + \dots, \quad n = 1, 2, 3, \dots$$

The inversion of a set of equations of this type is a familiar exercise in elementary number theory. The result is

$$(2.11) \quad r^n a_n = \mu_1 b_n + \mu_2 b_{2n} + \mu_3 b_{3n} + \dots, \quad n = 1, 2, 3, \dots,$$

the μ_n being the n th Möbius number defined by (see, for example, [4])

$$(2.12) \quad \begin{aligned} \mu_1 &= 1, \\ \mu_i &= (-1)^r \quad \text{if } i \text{ is a product of } r \text{ distinct prime numbers,} \\ \mu_i &= 0 \quad \text{otherwise.} \end{aligned}$$

The values of the first fifteen Möbius numbers are given here together with the values λ_k needed later on.

$$\begin{aligned} k &= 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, \\ \mu_k &= +1, -1, -1, 0, -1, +1, -1, 0, 0, +1, -1, 0, -1, +1, +1, \\ \lambda_k &= +1, +1, -1, 2, -1, -1, -1, 4, 0, -1, -1, -2, -1, -1, +1. \end{aligned}$$

If $f(x)$ is a real function of x , it follows that a_n is real. If we substitute into the real part of (2.11) the expression for b_n given by (2.9) and a_n by (2.1), we find (1.3), thus establishing Theorem 1.

3. Other theorems of a similar nature. There are several generalizations and modifications of the above derivation which lead to formulas of a nature similar to that of Theorem 1. In this section we outline these derivations and state corresponding results.

We introduce a set of rule sum operators in analogy to (1.1) above. These are “off set” trapezoidal rules,

$$(3.1) \quad R^{[n,\alpha]}g(t) = \frac{1}{n} \sum_{j=1}^n g\left(\frac{2j + \alpha - 1}{2n}\right), \quad |\alpha| < 1.$$

The Fourier series corresponding to one of these forms a generalization of the Poisson summation formula (see [2]), namely,

$$(3.2) \quad R^{[n,\alpha]}g(t) = \int_0^1 g(t) dt + \sum_{\substack{k=-\infty \\ k \neq 0}}^{\infty} e^{-\pi i(\alpha-1)k} \int_0^1 g(t) e^{2\pi i k n t} dt.$$

We are interested here in cases in which $\alpha = 0, \pm\frac{1}{2}, 1$. We define in analogy to (2.9) above,

$$\begin{aligned}
 b_n^{[0]} &= R^{[n,0]}f(re^{2\pi it}) - \int_0^1 f(re^{2\pi it}) dt, \\
 (3.3) \quad b_n^{[1]} &= b_n = R^{[n,1]}f(re^{2\pi it}) - \int_0^1 f(re^{2\pi it}) dt, \\
 b_n^{[2]} &= \frac{1}{2}(b_n^{[1]} - b_n^{[0]}), \\
 ib_n^{[3]} &= \frac{1}{2}(R^{[n,-1/2]}f(re^{2\pi it}) - R^{[n,1/2]}f(re^{2\pi it})).
 \end{aligned}$$

Substituting for a_n given by (2.3) or (2.4) into (3.3), we find respectively

$$\begin{aligned}
 (3.4) \quad b_n^{[0]} &= -r^n a_n + r^{2n} a_{2n} - r^{3n} a_{3n} + \dots, \\
 b_n^{[1]} &= r^n a_n + r^{2n} a_{2n} + r^{3n} a_{3n} + \dots, \\
 b_n^{[2]} &= r^n a_n + r^{3n} a_{3n} + r^{5n} a_{5n} + \dots, \\
 b_n^{[3]} &= r^n a_n - r^{3n} a_{3n} + r^{5n} a_{5n} - \dots, \quad n = 1, 2, 3, \dots.
 \end{aligned}$$

All these equations may be inverted in a similar way. The results are given in the following theorem.

THEOREM 2. *Under the hypotheses of Theorem 1,*

$$\begin{aligned}
 (3.5) \quad r^n a_n &= - \sum_{m=1}^{\infty} \lambda_m b_{mn}^{[0]}, \\
 r^n a_n &= \sum_{m=1}^{\infty} \mu_m b_{mn}^{[1]}, \\
 r^n a_n &= \sum_{m=1}^{\infty} \mu_{2m-1} b_{(2m-1)n}^{[2]}, \\
 r^n a_n &= i \sum_{m=1}^{\infty} (-1)^m \mu_{2m-1} b_{(2m-1)n}^{[3]},
 \end{aligned}$$

where the set of numbers λ_m are defined by

$$\begin{aligned}
 (3.6) \quad \lambda_m &= \mu_m, \quad m \text{ odd}, \\
 \lambda_m &= 2^{\alpha-1} \mu_{2n+1}, \quad m = 2^\alpha (2n + 1), \quad \alpha \geq 1.
 \end{aligned}$$

Any of the set (3.5) may be used to compute a_n . If $f(x)$ is a real function of x , the first three involve only the real part of $f(z)$ and the fourth only the imaginary part of $f(z)$. Whether or not $f(x)$ is a real function of x , the first two involve the function evaluation $f(0)$ while the second two do not require this function evaluation.

4. Convergence and degree. The rate of convergence of the series in Theorem 2 is indicated by the following theorem.

THEOREM 3. *For all p ,*

$$\lim_{m \rightarrow \infty} b_m^{[k]} m^p = 0, \quad k = 0, 1, 2, 3.$$

In other words the sequence $b_m^{[k]}$ approaches zero faster than any power of m . We present a proof for $k = 1$. The proof for the other values of k is similar.

Proof. The Euler Maclaurin summation formula states:

$$\begin{aligned} R^{[m,1]}h(t) - \int_0^1 h(t) dt &= 2 \sum_{s=1}^{p-1} (-1)^{s+1} \zeta(2s) \frac{h^{(2s-1)}(1) - h^{(2s-1)}(0)}{(2\pi m)^{2s}} \\ &+ \frac{2(-1)^p}{(2\pi m)^{2p}} \int_0^1 h^{(2p)}(t) \sum_{r=1}^{\infty} \frac{1 - \cos 2\pi rmt}{r^{2p}} dt, \end{aligned} \tag{4.1}$$

where $h(t)$ and all its derivatives are continuous in the interval $0 \leq t \leq 1$. Since $f(z)$ is analytic in a domain which includes the circle $|z| = r$, it follows that

$$h(t) = f(re^{2\pi it})$$

satisfies the above condition. In addition $h(t)$ is periodic and analytic. Thus,

$$h^{(2s-1)}(1) - h^{(2s-1)}(0) = 0 \tag{4.2}$$

and the $(2p)$ th derivative of $h(t)$ is bounded,

$$|h^{(2p)}(t)| \leq M_{2p}. \tag{4.3}$$

If we use the inequality

$$\left| \sum_{r=1}^{\infty} \frac{1 - \cos 2\pi rmt}{r^{2p}} \right| \leq 2\zeta(2p) \tag{4.4}$$

and the definition (3.3), we find without difficulty that

$$|b_m^{[1]}| \leq \frac{4M_{2p} \zeta(2p)}{(2\pi m)^{2p}}. \tag{4.5}$$

The theorem follows immediately from this.

The formulas of Theorem 2 express the n th derivative as an infinite sum of terms, each of which is a sum of function evaluations. In practice this infinite sum is approximated by the sum of the first M terms, and M is determined by the size of the terms and the apparent convergence rate of the series. The connection between M and the degree of the formula is indicated by the following theorem.

THEOREM 4. *The formula*

$$a_n = \frac{1}{r^n} \sum_{m=1}^M \mu_m b_m^{[1]} \tag{4.6}$$

and the corresponding finite sums for $k = 0, 2, 3$ are exact if $f(z)$ is a polynomial of degree d , where

$$(4.7) \quad \begin{aligned} d < (M + 1)n, & \quad k = 0, 1, \\ d < (2M + 1)n, & \quad k = 2, 3. \end{aligned}$$

Proof. If $f(x)$ is a polynomial of degree d , then

$$(4.8) \quad a_n = 0, \quad n > d,$$

and using (3.4),

$$(4.9) \quad b_n^{[k]} = 0, \quad n > d.$$

It follows that if d satisfies (4.7) the sum given in (4.6) is identical with the sum in (3.5) which, by Theorem 2, is exact.

5. A roundoff error estimate. From a practical standpoint, one of the most important features of these methods is that a simple estimate of the roundoff error is available at an early stage in the calculation. In this section we derive this estimate and indicate how it may be used to choose a value for the parameter r suitable to the required accuracy of the particular calculation.

We assume that the calculation is carried out using floating point arithmetic. We denote by ϵ the maximum relative roundoff error in the arithmetic operations and in the function evaluations. In the numerical examples below, the computer used 36-bit floating point significant arithmetic, and

$$(5.1) \quad \epsilon = 0.5 \times 10^{-10}.$$

As an illustration of these methods we consider the calculation of the fifth derivative $f^{(5)}(0)$ of the function

$$(5.2) \quad f(x) = \frac{e^x}{\sin^3 x + \cos^3 x}.$$

(See Table 1.) This example satisfies our requirements of being tedious to differentiate analytically, but easy to compute for complex arguments. It is clear that the value of the required derivative is integer. Using a symbolic algebraic processor [1]—certainly the most convenient way to solve this particular problem—we find

$$(5.3) \quad f^{(5)}(0) = -164.$$

Conventional methods of numerical differentiation using divided differences or methods using the Neville algorithm [3] are based implicitly on formulas of the type

$$(5.4) \quad f^{(n)}(0) = \sum a_i f(x_i),$$

TABLE I

Calculation of $f^{(5)}(0)$ where $f(x)$ is given by (5.2)

(a) The value of ϵ is given by (5.1).

(b) In each case, the value of δ_1 is available after the first three function evaluations.

(c) The series is terminated after the M th term, M being the first integer for which $|b_{Mn}| < \epsilon G_M$ and $\mu_M \neq 0$. In each case the sequence $|b_{Mn}|$ is strongly monotonic decreasing at this stage.

(d) The quantity ν is the number of times the subroutine for evaluating $f(z)$ was entered in our particular code. This is not the minimum possible, since we found it convenient to evaluate the function at a point each time this quantity was required. However, we did make use of the fact that $g(t)$ is symmetric about $t = \frac{1}{2}$ which approximately halves the number of function evaluations required.

| r | M | S_M | δ_1 | δ_M | $f^{(5)}(0)$ | $\delta_M f^{(5)}(0) $ | ν |
|-----|-----|--------------------------------|------------------------|------------------------|-----------------|-------------------------|-------|
| 0.1 | 3 | $-1.3666605809 \times 10^{-5}$ | 4.10×10^{-6} | 4.10×10^{-6} | -163.99926, 973 | 0.00067, 246 | 16 |
| 0.4 | 7 | $-1.3994666602 \times 10^{-2}$ | 6.52×10^{-9} | 6.34×10^{-9} | -163.99999, 925 | 0.00000, 104 | 59 |
| 0.7 | 37 | $-2.2969566684 \times 10^{-1}$ | 8.99×10^{-10} | 6.13×10^{-10} | -164.00000, 013 | 0.00000, 010 | 1097 |

in which the constants a_i vary considerably in magnitude and are of different signs. This causes small errors in the function values $f(x_i)$ to be amplified in a manner which is difficult to predict quantitatively. In the above example, a fifth degree polynomial through computed function values at six points evenly spaced in the interval $[-0.1, 0.1]$ gives

$$(5.5) \quad f^{(5)}(0) = -168.5 \dots$$

An exhaustive experiment using several dozen different point distributions and different intervals produced results both more and less accurate than this. The closest result was correct to five significant figures but some results were even of the wrong sign.

The methods described in this paper involve formulas of type (5.4), all the a_i having the same magnitude. Specifically the steps involved in using (1.3) are

$$(5.6) \quad R^{[mn,1]}g(t) = \frac{1}{mn} \sum_{j=1}^{mn} g(t_j), \quad t_j = \frac{j}{mn},$$

$$(5.7) \quad b_{mn} = R^{[mn,1]}g(t) - f(0),$$

$$(5.8) \quad S_M = \sum_{m=1}^M \mu_m b_{mn}, \quad |\mu_m| = 1 \text{ or } 0,$$

$$(5.9) \quad \frac{f^{(n)}(0)}{n!} = a_n \doteq \frac{S_M}{r^n}.$$

While each of these steps involves some roundoff error, it is clear that in (5.6), (5.7) and (5.8) there is no undue amplification of the absolute

value of the errors in $g(t_j)$. In fact, the averaging process in (5.6) may damp these errors down. Consequently, the computed S_M may be expected to have an absolute error of the order of ϵG_M , where $G_M = G_M(r)$ is the largest function value occurring in the calculation of S_M , that is,

$$(5.10) \quad G_M = \max \{ |f(0)|, |g(t_j)| \}, \quad t_j = j/mn, 0 \leq m \leq M.$$

The final step (5.9) retains the *relative* accuracy. Thus a guide to the relative accuracy in the final result is given by δ_M , where

$$(5.11) \quad \frac{\Delta f^{(n)}(0)}{f^{(n)}(0)} \doteq \delta_M = \frac{\epsilon G_M}{S_M} \doteq \frac{\epsilon G(r)}{r^n a_n},$$

where

$$(5.12) \quad G(r) = \max |g(t)| = \lim_{M \rightarrow \infty} G_M(r).$$

The possibly large size of this relative error is a consequence of (5.6) and (5.7). If $|f(0)|$ is small or zero, then (5.6) involves numbers $g(t_j)$ of different sign. If $|f(0)|$ is large, (5.7) may involve taking the difference of two nearly equal numbers. In either case, relative accuracy is lost, although absolute accuracy is retained.

The value δ_M is in no rigorous sense a bound. The actual accuracy may be less than δ_M if roundoff errors unexpectedly cancel. Or one may have a slightly larger error, particularly if M is large. For example, if $f(z)$ involves a large nearly constant component, then the accumulation of the sum in (5.6) may be an additional significant source of error. This could be avoided at little expense by using double precision accumulations.

The estimate δ_M is an a posteriori estimate. Remembering that δ_M is not required at all accurately, but that only its rough magnitude is needed, the term δ_1 is usually a sufficiently accurate approximation to δ_M and is therefore an estimate of the roundoff error in the final result.

The advantage in using δ_1 is that it provides an *early* indication of whether the chosen value of the parameter r is suitable in view of the required accuracy. If δ_1 indicates an unacceptably large error, the calculation may be recommenced at once using a larger value of r . Since the series converges more slowly if r is larger, the coding of the problem may be arranged to use ultimately as small a value of r as is convenient and also consistent with the acceptable roundoff error.

6. Conclusion. It seems inappropriate to compare in any detail the methods given here with standard methods requiring function evaluations with real argument only. From the point of view of the numerical analyst these methods show that, once complex arguments are allowed, the principal difficulties encountered in numerical differentiation simply disappear. Theo-

rem 4 indicates that it is trivial to invent methods of any degree and the discussion of §5 shows that the roundoff error may be controlled in a straightforward manner. The methods given here are not minimal in the sense relating to the degree of the polynomial approximation and for this reason may well be unacceptable to the more conventional numerical analyst.

For the practical computer user with a function to differentiate, these methods are quite attractive. They are based on simple formulas, are easy to code, converge rapidly and, perhaps most important, include a useful error estimate.

7. Acknowledgment. Both of the authors are happy to acknowledge help and support received during independent visits to the Institut für Angewandte Mathematik, Eidgenössische Technische Hochschule, Zürich, Switzerland, where most of this work was carried out.

REFERENCES

- [1] M. ENGELI, *Design and implementation of an algebraic processor*, Privately circulated report, Eidgenössische Technische Hochschule, Zürich, 1966.
- [2] J. N. LYNESS, *The calculation of Fourier coefficients*, this Journal, 4 (1967), pp. 301-315.
- [3] J. N. LYNESS AND C. B. MOLER, *Van der Monde systems and numerical differentiation*, Numer. Math., 8 (1966), pp. 458-464.
- [4] G. PÓLYA AND G. SZEGÖ, *Aufgaben und Lehrsätze aus der Analysis*, vol. 2, Springer-Verlag, Berlin, 1954.
- [5] E. C. TITCHMARSH, *Introduction to the Theory of Fourier Integrals*, Oxford University Press, Oxford, 1948.