In other words, $D(n)$ is a function of $n^2$ variables $a_{jk}$ which vary over the bounded and closed domain $\bar{D}: \{|\ a_{jk}\ | \leqq K\}$; hence this function is bounded and attains its maximum value on the boundary of the domain $\bar{D}$.

*Proof.* Let $a_{jk} = r_{jk}e^{i\theta_{jk}}$ and $A_{jk} = R_{jk}e^{i\phi_{jk}} =$ the co-factor of $a_{jk}$, where $K \geqq r_{jk} \geqq 0$ and $R_{jk} \geqq 0$. Then, expanding by the $j$th row, we have

$$
|\ D(n)\ | = \left|\ \sum_{k=1}^{n} a_{jk}A_{jk}\ \right| = \left|\ \sum_{k=1}^{n} r_{jk}R_{jk}\ e^{i(\theta_{jk}+\phi_{jk})}\ \right|
$$

(2)

$$
\leqq \sum_{k=1}^{n} r_{jk}R_{jk} \leqq \sum_{k=1}^{n} KR_{jk} = D'(n),
$$

where $D'(n)$ is the $n$th order determinant whose entries are

(3)
$$
a'_{jk} = \begin{cases} a_{jk}, & \text{if } r_{jk} = K \text{ and } \theta_{jk} + \phi_{jk} \equiv 0 \pmod{2\pi}, \\ Ke^{-i\phi_{jk}}, & \text{if } r_{jk} < K \text{ or } \theta_{jk} + \phi_{jk} \not\equiv 0 \pmod{2\pi}. \end{cases}
$$

By applying the same process to the other rows, we obtain a determinant $D^*(n)$ whose entries $|\ a^*_{jk}\ | = K$ and $|D^*(n)| \geqq |\ D(n)|$. Hence, $\text{Max}_{|a_{jk}|\leqq K}\ |\ D(n)| \leqq \text{Max}_{|a_{jk}|=K}\ |\ D(n)|$; thus the proof of the theorem can be completed since the reverse inequality is trivial.

50 Mohigan Drive
Oneonta, New York

# On the Numerical Solution of $y' = f(x, y)$ by a Class of Formulae Based on Rational Approximation

## By John D. Lambert and Brian Shaw

**1. Introduction.** Most finite difference formulae in common usage for the numerical solution of first-order differential equations are based on polynomial approximation. Two exceptions are the formulae based on exponential approximation proposed by Brock and Murray [1], and the formulae of Gautschi [2] which are derived from trigonometric polynomials. The use of rational functions as approximants has been studied by many authors, including Remes [3], Maehly [4] and Stoer [5], but the main concern of most of this work has been the direct approximation of a given function. Algorithms for interpolation based on rational functions have been proposed by Wynn [6], and methods for numerical integration and differentiation based on Padé approximation have been studied by Kopal [7]. It is the purpose of the present paper to derive a class of formulae, based on rational approximation, for the numerical solution of the initial value problem

(1)
$$
y' = f(x, y), \qquad y(x_0) = y_0 .
$$

The formulae proposed give exact results when the theoretical solution of (1) is a rational function of a certain degree, just as many of the classical difference formulae

give exact results when the theoretical solution is a polynomial of the appropriate degree.

Formulae based on polynomial approximation frequently give poor results if the integration of (1) is pursued too close to a singularity. It will be demonstrated that the class of formulae proposed can give better results in such circumstances.

**2. Derivation of the Formulae.** The method of derivation will first be used to obtain one of the well-known polynomial formulae.

Along the $x$-axis, consider the points $x_r$ to be given by

$$x_r = x_0 + rh \qquad\qquad (r = 0, 1, 2, \cdots),$$

where $h$ is the distance between consecutive points. The formula to be derived will predict values of $y_r$ which will approximate to $y(x_r)$, the theoretical solution of (1) at $x_r$. Let us assume that the solution of (1) is locally represented in the range $[x_n, x_{n+2}]$ by the polynomial

$$(2) \qquad\qquad F(x) = \sum_{s=0}^{4} a_s x^s.$$

This polynomial must pass through the points $(x_n, y_n)$, $(x_{n+1}, y_{n+1})$, $(x_{n+2}, y_{n+2})$, and, moreover, must assume at these points the slopes given by $y' = f(x, y)$. The following six equations must then be satisfied.

$$(3) \qquad\qquad F(x_{n+j}) = y_{n+j}, \qquad F'(x_{n+j}) = f_{n+j} \qquad\qquad (j = 0, 1, 2),$$

where $f_r = f(x_r, y_r)$. The eliminant of the five undetermined coefficients $a_s$ from the six equations (3) is the familiar Simpson's rule,

$$(4) \qquad\qquad y_{n+2} - y_n = \frac{h}{3}(f_{n+2} + 4f_{n+1} + f_n).$$

The same approximant (2) can also be used to derive a two-point formula involving higher derivatives of $f$, which can be calculated using (1). Thus

$$y^{(s+1)} = f^{(s)} \equiv \frac{d^s f}{dx^s} = (f^{(s-1)})_x + (f^{(s-1)})_y f \quad (s = 1, 2, 3, \cdots),$$

where $f^{(0)} \equiv f$. The eliminant of the $a_s$ from the six equations

$$F(x_{n+j}) = y_{n+j}, \qquad F'(x_{n+j}) = f_{n+j}, \qquad F''(x_{n+j}) = f^{(1)}_{n+j} \quad (j = 0, 1)$$

is

$$y_{n+1} - y_n = \frac{h}{2}(f_{n+1} + f_n) - \frac{h^2}{12}(f^{(1)}_{n+1} - f^{(1)}_n),$$

one of the class of formulae derived by Lambert and Mitchell [8].

This method of derivation, although tedious by comparison with that employed in [8], has the virtue that it can be applied when the basic polynomial approximant $F(x)$ is replaced by a rational approximant

$$R(x) = P(x)/Q(x),$$

where $P(x)$ is a polynomial of degree $p$, and $Q(x)$ a monic polynomial of degree $q$.

**3. Explicit Formulae.** Let $k + 1$ be the number of points utilised in the formula, and let $l$ be the order of the highest derivative of $y$ involved. Then an optimum class of explicit formulae can be derived if

(5) $$k(l + 1) + 1 = p + q + 2.$$

The most useful class of formulae is obtained by putting $k = q = 1$, that is, a class of two-point formulae based on a rational approximant whose denominator is a linear function. It follows from (5) that derivatives of up to order $p + 1$ must be employed. Applying the method described above to the approximant

(6) $$R(x) = \left( \sum_{s=0}^{p} a_s x^s \right) \bigg/ (b_0 + x),$$

the following eliminant is obtained.

(7) $$y_{n+1} - y_n = \sum_{s=1}^{p-1} \frac{h^s}{s!} f_n^{(s-1)} + \frac{h^p}{p!} \frac{(p + 1)(f_n^{(p-1)})^2}{(p + 1)f_n^{(p-1)} - hf_n^{(p)}}.$$

Taylor expansion of (7) shows that the associated principal truncation error is

$$\frac{h^{p+2}}{(p + 1)!} \left[ -\frac{y^{(p+2)}}{p + 2} + \frac{(y^{(p+1)})^2}{(p + 1)y^{(p)}} \right].$$

Each formula of class (7) is seen to be a Taylor series with a rational correcting term, and, being a two-point formula, cannot suffer strong instability.

The class of three-point formulae based on the approximant (6) will require, by (5), that $p = 2l$, that is, the numerator of the rational approximant must be of even degree. The general formula of this class is very unwieldy, and only the first two members are quoted.

$p = 2, q = 1:$ $\quad 3y_{n+2} - 4y_{n+1} + y_n = \dfrac{2h}{3}(2f_{n+1} + f_n)$

(8) $$+ \frac{4h^2}{3} \frac{(f_{n+1} - f_n)^2}{3(y_{n+1} - y_n) - h(f_{n+1} + 2f_n)},$$

Truncation error: $\quad h^4 \left( -\dfrac{1}{2} y^{(4)} + \dfrac{2}{3} \dfrac{(y^{(3)})^2}{y^{(2)}} \right);$

$p = 4, q = 1:$ $\quad y_{n+2} - y_n = \dfrac{2h}{9}(8f_{n+1} + f_n) + \dfrac{2h^2}{9}(2f_{n+1}^{(1)} - f_n^{(1)})$

(9) $$- \frac{4h^2}{9} \frac{[2(f_{n+1} - f_n) - h(f_{n+1}^{(1)} + f_n^{(1)})]^2}{18(y_{n+1} - y_n) - 2h(4f_{n+1} + 5f_n) + h^2(f_{n+1}^{(1)} - 2f_n^{(1)})},$$

Truncation error: $\quad h^6 \left( -\dfrac{1}{90} y^{(6)} + \dfrac{1}{75} \dfrac{(y^{(5)})^2}{y^{(4)}} \right).$

Expansion of the rational terms in formulae (8) and (9) shows these terms to be of orders $h^2$ and $h^4$, respectively. Thus formulae (8) and (9) have the interesting property that they do not exhibit strong instability in the Dahlquist sense. The corresponding formulae based on polynomial approximation, as derived in [8], are both strongly unstable.

For fixed values of $k$ and $l$, a formula based on an approximant with $p = p^*$, $q = q^*$ will give, for the problem (1), an algorithm identical with that which would be obtained from a formula based on an approximant with $p = q^*$, $q = p^*$, applied to the problem formed by applying the transformation $\bar{y} = 1/y$ to (1). It follows that the next approximant which should be considered will have $p = q = 2$. The simplest formula based on this approximant is given by $k = 1$, $l = 4$, and is

$$y_{n+1} - y_n = hf_n + h^2$$

$$(10) \qquad \cdot \frac{6f_n^{(1)}[(3f_n^{(1)})^2 - 2f_n f_n^{(2)}] + hf_n[3f_n^{(1)}f_n^{(3)} - 4(f_n^{(2)})^2]}{12[3(f_n^{(1)})^2 - 2f_n f_n^{(2)}] + 6h[f_n f_n^{(3)} - 2f_n^{(1)}f_n^{(2)}] + h^2[4(f_n^{(2)})^2 - 3f_n^{(1)}f_n^{(3)}]},$$

$$\text{Truncation error:} \qquad h^5\left[ -\frac{1}{120}y^{(5)} + \frac{1}{144}\frac{8(y^{(3)})^3 - 3y^{(1)}(y^{(4)})^2}{3(y^{(2)})^2 - 2y^{(1)}y^{(3)}} \right].$$

It should be noticed in passing that formula (10) and the formulae of class (7) are applicable as quadrature formulae for the evaluation of $\int_a^b f(x)\,dx$, particularly when the range of integration is close to, or includes, a singularity of $f(x)$.

If, during a calculation, the denominator of the rational term in any of the proposed formulae becomes zero (or very nearly so) at a station at which it is known that the theoretical solution of (1) does not have a singularity, then the mesh length must be altered, or, if this remedy fails, a different formula must be used. (In this context, it would appear that two-point formulae cause less trouble than do those with three points.) A change of sign of the denominator would indicate that a pole of the approximant $R(x)$ had fallen within the local range of application of the formula—a situation analogous to that of a polynomial approximant which becomes oscillatory within the local range of application of the associated formula. Although it can be argued that the local intervention of a pole is potentially more serious than the onset of an oscillation, the formulae based on rational approximation have the advantage that the occurrence of this difficulty during a calculation is easily detected by keeping a separate check on the behaviour of the denominator. The onset of polynomial oscillation in classical formulae is much more difficult to detect.

**4. Implicit Formulae.** Implicit formulae can be obtained if

$$(11) \qquad (k + 1)(l + 1) = p + q + 2.$$

Two-point formulae with $q = 1$ can therefore be obtained only if $p$ is odd, and then derivatives of $y$ up to order $(p + 1)/2$ must be employed. The first two formulae in this class are

$$p = 1, q = 1: \qquad y_{n+1} - y_n = h^2 \frac{f_n f_{n+1}}{y_{n+1} - y_n},$$

$$(12)$$

$$\text{Truncation error:} \qquad h^3\left[ -\frac{y^{(3)}}{6} + \frac{(y^{(2)})^2}{4y^{(1)}} \right];$$

$$p = 3, q = 1: \qquad y_{n+1} - y_n$$

$$(13) \qquad = -h^2 \frac{4(f_{n+1} - f_n)^2 + 12f_n f_{n+1} + 2h(f_n f_{n+1}^{(1)} - f_n^{(1)} f_{n+1}) + h^2 f_n^{(1)} f_{n+1}^{(1)}}{12(y_{n+1} - y_n) - 12h(f_{n+1} + f_n) - 2h^2(f_{n+1}^{(1)} - f_n^{(1)})},$$

$$\text{Truncation error:} \qquad h^5\left[ -\frac{1}{720}y^{(5)} + \frac{1}{576}\frac{(y^{(4)})^2}{y^{(3)}} \right].$$

It is of interest to observe that (12) equates $(y_{n+1} - y_n)/h$ to the geometric mean of $f_n$ and $f_{n+1}$, while the corresponding formula based on a polynomial approximant equates the same expression to the arithmetic mean of $f_n$ and $f_{n+1}$.

In view of the strong stability of the three-point explicit formulae already derived, there would appear to be little point in quoting comparable implicit formulae, which, in any case, turn out to be excessively unwieldy.

**5. Numerical Results.** The example used to illustrate the formulae derived above is the initial value problem

$$(14) \qquad\qquad y' = 1 + y^2, \qquad y(0) = 1,$$

whose theoretical solution is $y = \tan(x + \pi/4)$.

A comparison is made on the basis of two-point formulae. Problem (14) is solved, first using formula (10), and secondly using the formula of class (7) obtained by setting $p = 3$. This gives

$$(15) \qquad y_{n+1} - y_n = hf_n + \frac{h^2}{2} f_n^{(1)} + \frac{h^3}{6} \frac{4(f_n^{(2)})^2}{4f_n^{(2)} - hf_n^{(3)}},$$

$$\text{Truncation error:} \qquad \frac{h^5}{24}\left[ -\frac{1}{5} y^{(5)} + \frac{1}{4} \frac{(y^{(4)})^2}{y^{(2)}} \right].$$

Both of these formulae are explicit, involve two points, and utilise derivatives of $y$ up to order four. The corresponding formula based on a polynomial approximant is the truncated Taylor series formula,

$$(16) \qquad y_{n+1} - y_n = hf_n + \frac{h^2}{2} f_n^{(1)} + \frac{h^3}{6} f_n^{(2)} + \frac{h^4}{24} f_n^{(3)},$$

$$\text{Truncation error:} \qquad -\frac{1}{120} h^5 y^{(5)}.$$

TABLE I
*Two-point Formulae*

| $x$ | Theoretical Solution | Formula (16) Polynomial | Formula (15) Rational $p = 3$, $q = 1$ | Formula (10) Rational $p = 2$, $q = 2$ |
|---|---|---|---|---|
| 0 | 1.000,000,000 | 1.000,000,000 | 1.000,000,000 | 1.000,000,000 |
| 0.05 | 1.105,355,590 | 1.105,354,167 | 1.105,355,556 | 1.105,355,575 |
| 0.10 | 1.223,048,880 | 1.223,045,160 | 1.223,048,805 | 1.223,048,846 |
| 0.15 | 1.356,087,851 | 1.356,080,366 | 1.356,087,728 | 1.356,087,792 |
| 0.20 | 1.508,497,647 | 1.508,483,855 | 1.508,497,464 | 1.508,497,556 |
| 0.25 | 1.685,796,417 | 1.685,771,749 | 1.685,796,159 | 1.685,796,284 |
| 0.30 | 1.895,765,123 | 1.895,720,992 | 1.895,764,765 | 1.895,764,932 |
| 0.35 | 2.149,747,640 | 2.149,667,006 | 2.149,747,147 | 2.149,747,367 |
| 0.40 | 2.464,962,757 | 2.464,809,445 | 2.464,962,070 | 2.464,962,364 |
| 0.45 | 2.868,884,028 | 2.868,574,494 | 2.868,883,051 | 2.868,883,451 |
| 0.50 | 3.408,223,442 | 3.407,542,560 | 3.408,222,003 | 3.408,222,567 |
| 0.55 | 4.169,364,046 | 4.167,671,633 | 4.169,361,803 | 4.169,362,642 |
| 0.60 | 5.331,855,223 | 5.326,819,985 | 5.331,851,409 | 5.331,852,773 |
| 0.65 | 7.340,436,575 | 7.320,574,452 | 7.340,429,058 | 7.340,431,623 |
| 0.70 | 11.681,373,800 | 11.552,695,821 | 11.681,353,989 | 11.681,360,445 |
| 0.75 | 28.238,252,850 | 25.710,677,828 | 28.238,132,170 | 28.238,169,733 |

Formulae (10), (15) and (16) all have the same order of principal truncation error. The numerical solutions of (14) by these three formulae are given in Table I, together with the theoretical solution. A mesh length of $h = .05$ was used, allowing the last station to be $x = .75$, whereas the singularity of the theoretical solution is at $x = .7854$. The calculations were done on an IBM 1620 computer, working, in floating point, to fourteen decimal places.

It is seen from Table I that the performance of the formulae based on rational approximants is markedly better than that of the polynomial-based formula.

The same problem is solved again by the three-point formula (9), in order to illustrate a remark made in Section 3. The corresponding optimum formula based on a polynomial interpolant, using the same points and derivatives as (9), is shown in [8] to be strongly unstable. The nearest stable polynomial formula for comparison purposes is the following, also taken from [8].

$$y_{n+2} - y_n = 2hf_n + \frac{2h^2}{3}(2f_{n+1}^{(1)} + f_n^{(1)}),$$

(17)

$$\text{Truncation error:} \quad -\frac{2}{45}h^5 y^{(5)}.$$

Table II shows the numerical solutions of (14) by these two formulae. This time the polynomial formula is better than the rational, but neither is good. However, a separate print-out of the denominator of the rational term in (9), quoted in the last column of Table II, shows that the pole of the rational approximant intervenes frequently, indicating that formula (9) is unsuitable for the problem in hand. (Indeed, it can be seen that formula (9) gives a better result than (17) up to $x = .20$, where the pole of the approximant intervenes for the first time.) The denominators

TABLE II
*Three-point Formulae*

| $x$ | Theoretical Solution | Formula (17) Polynomial | Formula (9) Rational $p = 4$, $q = 1$ | Denominator |
|---|---|---|---|---|
| 0 | 1.000,000 | 1.000,000* | 1.000,000* | |
| 0.05 | 1.105,356 | 1.105,356* | 1.105,356* | |
| 0.10 | 1.223,049 | 1.223,039 | 1.223,049 | −0.000,044,308 |
| 0.15 | 1.356,088 | 1.356,073 | 1.356,088 | −0.000,061,787 |
| 0.20 | 1.508,498 | 1.508,462 | 1.508,506 | −0.000,084,800 |
| 0.25 | 1.685,796 | 1.685,738 | 1.681,962 | +0.000,010,870 |
| 0.30 | 1.895,765 | 1.895,652 | 1.894,095 | −0.064,401,230 |
| 0.35 | 2.149,748 | 2.149,548 | 2.144,446 | +0.047,904,767 |
| 0.40 | 2.464,963 | 2.464,571 | 2.459,798 | −0.053,803,086 |
| 0.45 | 2.868,884 | 2.868,107 | 2.858,529 | +0.023,417,882 |
| 0.50 | 3.408,223 | 3.406,506 | 3.392,591 | −0.058,975,769 |
| 0.55 | 4.169,364 | 4.165,158 | 4.139,087 | −0.026,305,965 |
| 0.60 | 5.331,855 | 5.319,554 | 5.265,304 | −0.119,863,183 |
| 0.65 | 7.340,437 | 7.293,760 | 7.154,805 | −0.277,808,720 |
| 0.70 | 11.681,374 | 11.404,247 | 10.924,394 | −0.872,469,219 |
| 0.75 | 28.238,253 | 23.995,397 | 21.269,964 | −3.540,801,901 |

Starting values are marked with an asterisk.

of the rational terms of the two-point formulae (10) and (15), on the other hand, are of constant sign throughout the computation.

**Acknowledgment.** Mr. Shaw's share of the work of this paper was done whilst he was in receipt of a Carnegie Scholarship.

Department of Applied Mathematics
University of St. Andrews
St. Andrews, Scotland

1. P. BROCK & F. J. MURRAY, "The use of exponential sums in step by step integration," MTAC, v. 6, 1952, pp. 63–78. MR **13**, 873.
2. W. GAUTSCHI, "Numerical integration of ordinary differential equations based on trigonometric polynomials," *Numer Math.*, v. 3, 1961, pp. 381–397. MR **25** #1647.
3. E. REMES, "Sur le calcul effectif des polynomes d'approximation de Tchebycheff," *C. R. Acad. Sci. Paris*, v. 199, 1934, pp. 337–340.
4. H. MAEHLY, "Methods for fitting rational approximations," Parts II & III, *J. Assoc. Comput. Mach.*, v. 10, 1963, pp. 257–277. MR **28** #707.
5. J. STOER, "A direct method for Chebyshev approximation by rational functions," *J. Assoc. Comput. Mach.*, v. 11, 1964, pp. 59–69.
6. P. WYNN, "Über einen Interpolations-Algorithmus und gewisse andere Formeln, die in der Theorie der Interpolation durch rationale Funktionen bestehen," *Numer. Math.*, v. 2, 1960, pp. 151–182. MR **23** #B1636.
7. Z. KOPAL, "Operational methods in numerical analysis based on rational approximation," *Proceedings of a Symposium on Numerical Approximation* (1958), R. E. LANGER (ED.), Publication No. 1, Mathematics Research Center, Univ. of Wisconsin, Univ. of Wisconsin Press, Madison, 1959, pp. 25–43.
8. J. D. LAMBERT & A. R. MITCHELL, "On the solution of $y' = f(x, y)$ by a class of high accuracy difference formulae of low order," *Z. Angew. Math. Phys.*, v. 13, 1962, pp. 223–232. MR **25** #3610.

# The Numerical Solution of Eigenvalue Problems

## By Theodore R. Goodman

**1. Introduction.** One method for solving eigenvalue problems on a digital computer is to convert the governing differential equations to finite difference equations, apply the boundary conditions at either end of the interval, and form a secular equation for the unknown parameter (the eigenvalue) by setting the determinant associated with the resulting set of homogeneous algebraic equations for the ordinates of the solution equal to zero. Another way of solving eigenvalue problems is to use the Galerkin method. This consists of assuming the solution to be expanded in a complete set of functions satisfying the boundary conditions; upon introducing the series into the differential equation and requiring the error to be orthogonal to the functions in the set there results an infinite set of homogeneous equations for the coefficients. The secular equation is formed by setting the associated determinant equal to zero. These formulations invariably require the determination of the roots of a determinant of large order. The methods arise naturally out of the very nature of an eigenvalue problem and are seen to utilize the capability of digital computers to manipulate matrices of large order.

A completely different method for solving eigenvalue problems will be presented