# Notes on the Adaptive Simpson Quadrature Routine

J. N. LYNESS

*Argonne National Laboratory, Argonne, Illinois*

ABSTRACT. The adaptive Simpson quadrature routine is described, and the properties of the approximation error are investigated on an elementary level. Even this superficial investigation indicates some changes which might improve the efficiency of the routine. Four specific modifications in the routine are suggested to implement these changes.

KEY WORDS AND PHRASES: numerical integration, automatic integration, adaptive Simpson routine

CR CATEGORIES: 5.1, 5.11, 5.16

*Introduction*

A definitive account of the principles of automatic integration appears in the book *Numerical Integration* by Davis and Rabinowitz [2]. Several algorithms and routines are described which can be used to obtain an approximation having a prescribed tolerance $\epsilon$ to a given definite integral. In general, the user provides only the limits of integration $A$ and $B$, the desired tolerance $\epsilon$, and a subroutine $FUN(X)$ which calculates the integrand. The automatic integration scheme, in the form of a subroutine, provides a result, thus "relieving the user of any need to think."

One of the earlier, and certainly one of the most successful, routines of this type is the adaptive Simpson routine, published in algorithm form in 1962 (McKeeman [6]). Modifications of this routine (McKeeman [7], [9] and McKeeman and Tesler [8]) have been widely used in computing centers ever since. However, to the author's knowledge there is no account of this routine or error analysis of it in the open literature. In this paper the method on which the routine is based is investigated with a view to making improvements.

It is perhaps pertinent to mention the form improvements in an automatic integration routine might take. In general, the required tolerance or accuracy $\epsilon$ is given and a "good" routine attempts to attain only this accuracy using as few function evaluations as possible. In practice, the routine produces a result of greater accuracy. It is a defect if greater accuracy is obtained at a cost of additional function evaluations, even if the increase in accuracy is considerable and the number of additional function evaluations is marginal. In fact one modification mentioned above was introduced by McKeeman in an attempt to correct to some extent a defect of this type. Bearing this in mind, improvements result if any changes of the following nature can be made:

(i) a change which produces a result of lower accuracy, but still within the prescribed accuracy $\epsilon$, with a consequent reduction in the number of function evaluations required;

(ii) a change which produces a result of greater accuracy but with no increase in the number of function evaluations;

(iii) a change which produces a result of the same accuracy with a reduction in the number of function evaluations required.

It is clear that a useful way to investigate the possibility of making changes of this type is to obtain expressions for the discretization error; this is done in Section 3. Prior to this a brief description of the adaptive Simpson quadrature routine (ASQ) in its original form is given (Section 1). This is necessary to provide proper background as well as a consistent mathematical notation for use in the subsequent error analysis. In Section 3 two independent modifications of the routine are suggested which produce changes of type (i) above. One of these (modification 2) is simply to base the routine on interval bisection rather than trisection. In Section 4 a modification of type (ii), which consists of using information already available to produce a fifth order result rather than a third order result, is suggested. In Section 5 a type (iii) modification (which is essential if round-off error in the function values is significant but void otherwise) is suggested.

## 1. The Original Adaptive Simpson Quadrature Scheme (McKeeman 1962)

In its original form this scheme consists of a method to evaluate an approximation $R_{AS}f$ to

$$If = I[A, A + H]f(x) = \int_{A}^{A+H} f(x)dx$$

having an $L_1$ accuracy $\epsilon$; i.e. the approximation error $| R_{AS}f - If |$ is not to exceed

$$\epsilon \int_{A}^{A+H} | f(x) | dx. \tag{1.1}$$

To avoid introducing unnecessary complications, we describe the scheme in the trivially modified form: the absolute magnitude of the error should not exceed $\epsilon$; i.e.

$$| R_{AS}f - If | < \epsilon. \tag{1.2}$$

This scheme is based on Simpson's three-point and seven-point quadrature rules applied to subintervals of the interval $[A, A + H]$. We define these rules in operator notation:

$$R[a, a + h]f(x) = \frac{h}{6} \{f(a) + 4f(a + h/2) + f(a + h)\}, \tag{1.3}$$

$$R^{(3)}[a, a + h]f(x) = R[a, a + h/3]f(x) + R[a + h/3, a + 2h/3]f(x) \\ + R[a + 2h/3, a + h]f(x). \tag{1.4}$$

The operator $R^{(m)}[a, a + h]f(x)$ may be analogously defined as the result of dividing the interval $[a, a + h]$ into $m$ equal subintervals and applying the same rule $R$ to each. All these expressions are approximations to $I[a, a + h]f(x)$ such that the higher the value of $m$ the better (hopefully) is the approximation and the larger is the number $(2m + 1)$ of function values required to construct this approximation. The ASQ routine produces a result of the following form. We denote the interval

$[A, A + H]$ as the unique level 0 interval. This is trisected into three equal intervals $[A, A + H/3]$, $[A + H/3, A + 2H/3]$, and $[A + 2H/3, A + H]$ and these are termed level 1 intervals having length $H/3$. Some, but not necessarily all, of these level 1 intervals may be trisected into level 2 intervals having length $H/9$. Proceeding in this manner, the interval $[A, A + H]$ is ultimately subdivided into $n$ subintervals $[a_i, a_{i+1}]$, $i = 0, 1, \cdots, n - 1$, which span the complete interval. Thus

$$A = a_0 < a_1 < a_2 \cdots < a_{n-1} < a_n = B = A + H. \qquad (1.5)$$

The subintervals are in general of different lengths. A level $r$ subinterval has length

$$a_i - a_{i-1} = h_i = H3^{-r}. \qquad (1.6)$$

All the routines described subsequently use this type of subdivision or an analogous type based on bisection. They differ in that different "practical convergence criteria" lead to different subdivisions of the same type.

The result returned by the routine is

$$R_{As}[A, B]f(x) = \sum_{i=1}^{n} R^{(3)}[a_{i-1}, a_i]f(x), \qquad (1.7)$$

the sum of the approximations resulting from applying Simpson's seven-point rule to each subinterval in turn. This result requires $6n + 1$ function evaluations, and the subdivision method described above requires $n$ to be odd.

The actual subdivision, eq. (1.5), which depends on the function values encountered and the application of a criterion based on these values, is the *adaptive* part of the routine. The general intention is to obtain small intervals (high-level numbers) where the function is varying rapidly and large intervals where the function varies only slightly. The principal difference between the ASQ routine and the various modified versions described below is the method by which this end is attained. In the ASQ routine a set of level error constants

$$\epsilon(r) = \epsilon/3^r \qquad (1.8)$$

is assigned to each level, $\epsilon$ being the total error (see (1.2)). At a particular stage in the calculation, the interval $[a_i, a_i + h]$ is being considered. At this stage, the numbers $a_0, a_1, \cdots, a_{i-1}$ have been calculated and the approximation

$$\sum_{j=1}^{i} R^{(3)}[a_{j-1}, a_j]f(x) \simeq \int_A^{a_i} f(x) \, dx \qquad (1.9)$$

is available. The rest of the values for the ultimate subdivision, i.e. the value of $n$ and the numbers $a_{i+1}, a_{i+2}, \cdots, a_n$, have not been calculated. (However, a subset of these numbers, together with corresponding function values, is available; thus an increasingly good approximation to $\int_A^{A+H} |f(x)| \, dx$, but less accurate than the ultimate result, is available for use in the test in (1.1) if this test is to be used instead of (1.2).)

The expression $R[a_i, a_i + h]f(x)$ is usually available from a previous part of the calculation, together with function values $f(a_i), f(a_i + h/2), f(a_i + h)$. The expression $R^{(3)}[a_i, a_i + h]f(x)$ is now calculated, requiring four additional function evaluations. The magnitude of the difference between these two approximations to $I[a_i, a_i + h]f(x)$ is compared with $\epsilon(r)$. If this is less than $\epsilon(r)$, that is, if

$$R[a_i, a_i + h]f(x) - R^{(3)}[a_i, a_i + h]f(x) = \theta\epsilon(r), \qquad |\theta| \leq 1, \qquad (1.10)$$

then we allow that "this interval has converged." We define $h_{i+1} = h$, add $R^{(3)}[a_i, a_i + h]f(x)$ to the sum (1.9), and go on to consider the next interval to the right of $a_{i+1}$ where

$$a_{i+1} = a_i + h_{i+1}. \tag{1.11}$$

On the other hand, if the difference between these two approximations is greater than $\epsilon(r)$, i.e. $|\theta| > 1$ in (1.10), we proceed to subdivide the interval $[a_i, a_i + h]$ into three intervals $[a_i, a_i + h/3]$, $[a_i + h/3, a_i + 2h/3]$, $[a_i + 2h/3, a_i + h]$. Since the interval $[a_i, a_i + h]$ is a level $r$ interval, each of these subintervals is a level $r + 1$ interval. Also, in the calculation of $R^{(3)}[a_i, a_i + h]$ we have calculated $R[a_i, a_i + h/3]f(x)$ and the corresponding quantities for the other two subintervals. This latter information is stored for future use and we proceed to treat the interval $[a_i, a_i + h/3]$ in the manner just described for the interval $[a_i, a_i + h]$.

The only exception to this procedure is in dealing with the initial interval $[A, A + H]$. Here it is pretended that the "level 0 interval does not converge." Thus the routine is always required to base its result on at least 19 function values.

It should be noted that following this procedure carries out the integration from left to right. When the level $r$ interval $[a_i, a_i + h]$ is being treated, the approximation to the previous part of the integral (1.9) is already available, and information corresponding to 0, 1, or 2 intervals at each of the levels $r, r - 1, \cdots, 1$, all these intervals being to the right of $a_i + h$, is stored for subsequent attention.

The routine concludes when an interval $[a_{n-1}, B]$ converges. Since each interval $[a_i, a_{i+1}]$ has converged, we may write

$$R[a_i, a_{i+1}]f(x) - R^{(3)}[a_i, a_{i+1}]f(x) = \theta_{i+1}\epsilon(r) \tag{1.12}$$

where $|\theta_{i+1}| < 1$. Since

$$h_{i+1} = a_{i+1} - a_i = H3^{-r} \tag{1.13}$$

we have

$$\epsilon(r) = \epsilon(a_{i+1} - a_i)/H. \tag{1.14}$$

Thus

$$\sum_{i=0}^{n-1} R[a_i, a_{i+1}]f(x) - \sum_{i=0}^{n-1} R^{(3)}[a_i, a_{i+1}]f(x) = \sum_{i=0}^{n-1} \frac{\epsilon}{H}(a_{i+1} - a_i)\theta_{i+1}, \tag{1.15}$$

and the magnitude of the quantity on the right may be bounded by

$$\left| \sum_{i=0}^{n-1} \frac{\epsilon}{H}(a_{i+1} - a_i)\theta_{i+1} \right| \le \sum_{i=0}^{n-1} \frac{\epsilon}{H}(a_{i+1} - a_i)|\theta_{i+1}| \le \sum_{i=0}^{n-1} \frac{\epsilon}{H}(a_{i+1} - a_i) = \epsilon. \tag{1.16}$$

The second expression on the left-hand side of (1.15) is the result (1.7) returned by the routine based on $6n + 1$ function values. This evidently differs from a similar result based on $2n + 1$ function values by a quantity in magnitude less than $\epsilon$. In the absence of any evidence to the contrary, from a practical standpoint there is a strong presumption that the true integral differs from the result based on $6n + 1$ function values by less than $\epsilon$.

## 2. The Modified Adaptive Simpson Quadrature Scheme (McKeeman 1963)

The scheme described in detail in Section 1 was found in practice to be "over cautious"; that is, an accuracy considerably greater than $\epsilon$ was usually obtained. The scheme was modified by replacing the error level constants $\epsilon(r)$ given by (1.8) by

$$\epsilon(r) = \epsilon/(\sqrt{3})^r \tag{2.1}$$

but otherwise leaving the routine as it stands. This modification (which now appears to be standard) evidently gives better results [7].

The description in Section 1 applies equally to this routine, except that equations (1.14)–(1.16) are no longer applicable. It is not correct therefore to assert that the result requiring $6n + 1$ function values differs from a similar result requiring $2n + 1$ function values by less than $\epsilon$.

## 3. Modifications 1 and 2

In this section we carry out an elementary error analysis for ASQ-type quadrature routines. The purpose of this analysis is to justify certain modifications for these routines. No error bounds are given.

The theory presented here is valid for all integrand functions $f(x)$ which together with their first five derivatives are continuous in the interval of integration, to wit, $f(x) \in C^{(n)}[A, B]$, $n = 5$. If $n > 5$, various terms stated to be $O(h^5)$ are in fact $O(h^6)$. A trivially modified form of this theory is also valid if $n = 4$, but $f^{(4)}(x)$ is Lipschitz continuous in $[A, B]$ and the conclusions apply in this case also.

It should be noted that success in integrating through the singularities of functions which become infinite has been reported. An example is $\int_{-A}^{A} |\sin(1/x)|^{-1/2} dx$. The theory given here does not shed any light on the question, Why does the ASQ routine work for such functions? The modifications suggested below may turn out to hinder rather than help the routine in cases such as these.

In cases in which $f(x)$ is finite but its early derivatives have a finite number of discontinuities at known points, a standard procedure is to divide the integration interval at these points into subintervals and deal with each separately. If these points coincide with points used for function evaluation at an early level, and if $f(x)$ is continuous, there is no need to do even this. The local nature of the routine should ensure that this subdivision is carried out automatically.

A useful tool for dealing with quadrature discretization error is the Euler-Maclaurin summation formula (see (A.2)). In its conventional form this refers only to the trapezoidal rule. However, a similar formula exists for any rule (see Lyness and Ninham [5]). In particular, a special case ($q = 5$) of the appropriate formula (for Simpson's rule) may be stated as follows. If

$$f(x) = C^{(n)}[a, a + h], \qquad n \geq 5,$$

then

$$R^{(m)}[a, a + h]f(x) - \int_a^{a+h} f(x) \, dx = \frac{c_4 h^4}{m^4} \int_a^{a+h} f^{(4)}(x) \, dx \tag{3.1}$$

$$+ \frac{h^5}{m^5} \int_a^{a+h} g_5\left(m \frac{x - a}{h}\right) f^{(5)}(x) \, dx.$$

A brief derivation of this, together with expressions for $c_4$ and $g_5(x)$, are given in the Appendix. It follows trivially from (3.1) that

$$R[a, a + h]f(x) - R^{(m)}[a, a + h]f(x) = \frac{c_4 h^4 (m^4 - 1)}{m^4} \int_a^{a+h} f^{(4)}(x) \, dx$$

$$+ \frac{h^5}{m^5} \int_a^{a+h} \left\{ m^5 g_5 \left( \frac{x - a}{h} \right) - g_5 \left( m \frac{x - a}{h} \right) \right\} f^{(5)}(x) \, dx. \tag{3.2}$$

The left-hand side of (3.1) is the discretization error which we would like to bound. The left-hand side of (3.2) is the difference between two successive approximations, which can be calculated and whose value is used as a convergence criterion. The leading terms of both (3.1) and (3.2) are $O(h^5)$ and stand in the ratio $1 : (m^4 - 1)$.

A rough idea of the magnitudes of the various quantities involved at the points in the adaptive Simpson routines when the convergence criterion is applied may be obtained if in these formulas the terms of order $h^6$ are disregarded. In the version described in Sections 1 and 2, by comparing two cases of eq. (3.1), i.e. $m = 1$ and $m = 3$, we see that the seven-point approximation gives a result about 81 times as accurate as the three-point approximation. And from (3.2) we see that the difference between these approximations is about 80 times as large as the error in the seven-point approximation. Thus the requirement (eq. (1.10)) that for each convergent interval $| Rf - R^{(3)}f | < \epsilon(r)$ and the requirement $| R^{(3)}f - If | < \epsilon(r)/80$ are roughly equivalent. The net result, when account of all the convergent intervals is taken, is that an accuracy $\epsilon/80$ rather than $\epsilon$ is attained at a cost of approximately three times the number of function evaluations necessary to obtain the required accuracy $\epsilon$. Once this over-cautious property of the routine is recognized, it may be removed by a simple modification.

Modification 1.   Replace (1.8) by

$$\epsilon(r) = 80 \, \epsilon/3^r = (m^4 - 1)\epsilon/m^r.$$

A second modification is to replace $m = 3$ by $m = 2$. In the versions of Sections 1 and 2, if convergence in an interval is not attained, the routine triples the number of function evaluations to obtain a result 81 times as accurate. If $m = 2$, the routine would instead double the number of function evaluations to obtain a result 16 times as accurate. This seems to the author to be obviously preferable in intervals in which the convergence criterion fails to be satisfied only by a small margin. In the other cases, it does not seem to make any predictable difference.

Modification 2.   Replace $m = 3$ by $m = 2$ throughout the routine.

Incidentally, an adaptive method based on bisection rather than trisection is easier to implement as an automatic code.

A routine based on the adaptive Simpson routine but incorporating these particular modifications may be theoretically described in terms almost identical to those in Section 1. The essential difference is that (1.6), (1.8), (1.13), and (1.14) are replaced by

$$a_i - a_{i-1} = h_i = H2^{-r}, \qquad \epsilon(r) = 15\epsilon 2^{-r}. \tag{3.3}$$

This leads to the definition

$$R_{AS}[A, A + H]f(x) = \sum_{i=0}^{n-1} R^{(2)}[a_i, a_{i+1}]f(x). \tag{3.4}$$

Here

$$R_{AS}[A, A + H]f(x) - \int_A^{A+H} f(x)\, dx = \frac{1}{15} \sum_{i=0}^{n-1} \theta_{i+1}\epsilon(r)$$

$$+ \sum_{i=1}^{n} \frac{h_i^5}{15} \int_{a_{i-1}}^{a_i} \left\{ \frac{1}{2} g_5 \left( 2 \frac{x - a_{i-1}}{h_i} \right) - g_5 \left( \frac{x - a_{i-1}}{h_i} \right) \right\} f^{(5)}(x)\, dx \quad (3.5)$$

where

$$\theta_{i+1}\epsilon(r) = R[a_i, a_{i+1}]f(x) - R^{(2)}[a_i, a_{i+1}]f(x) \quad (3.6)$$

and

$$|\theta_i| < 1. \quad (3.7)$$

In conclusion we note the following.

(i) The approximation $R_{AS}[A, A + H]f(x)$ given by (3.4) based on $4n + 1$ function values differs from a similar approximation based on $2n + 1$ function values by less than $15\epsilon$ (not $\epsilon$, as in Section 1).

(ii) Since $R_{AS}[A, A + H]f(x)$ is approximately 16 times as accurate as this similar approximation, it is hoped that $R_{AS}[A, A + H]f(x)$ differs from the exact integral by less than $\epsilon$. This result is not guaranteed, as of course the terms of order $h^6$ in (3.5) may take any value.

(iii) The theory given here applies to integrands $f(x) \in C^{(5)}[A, A + H]$, assumes no upper limit on the possible level $r$, and disregards the effect of round-off error.

## 4. A Fifth-Order Approximation

As long as $f(x) \in C^{(5)}[A, A + H]$, the approximation $R_{AS}[A, A + H]f(x)$ has a discretization error given by equation (3.5). This error consists of the term

$$A = \frac{1}{15} \sum \theta_{i+1}\, \epsilon(r) \sim O\left( \sum h_i^5 \right), \quad (4.1)$$

which is smaller in magnitude than $\epsilon$, together with a term of order $\sum h_i^6$. However, the actual values of the individual terms $\theta_{i+1}\, \epsilon(r)$ in (4.1) are known, as the calculation of the right-hand side of (3.6) has been carried out to see if the interval has converged.

An insignificant amount of additional effort is required to make the actual value of term (4.1) available at the end of the calculation. Consequently, instead of returning a result containing an error $A + O\left( \sum h_i^6 \right)$ where $|A| < \epsilon$, since the actual value of $A$ may be precisely calculated, it seems more reasonable to adjust the result appropriately and return a result whose error is simply $O\left( \sum h_i^6 \right)$.

Modification 3.   The routine should return the result

$$Q_{AS}[A, A + H]f(x) = R_{AS}[A, A + H]f(x) - \frac{1}{15} \sum_{i=0}^{n-1} \theta_{i+1}\epsilon(r), \quad (4.2)$$

and also return the quantity $(1/15) \sum_{i=0}^{n-1} \theta_{i+1}\epsilon(r)$ as an additional error estimate.

It is quite straightforward to show that the result $Q_{AS}f$ has several conventional and more familiar interpretations. It is in fact the result obtained by applying the Newton-Cotes five-point formula to each interval. Thus, defining

$$Q[a, a + h]f(x) = \frac{2h}{45} \{7f(a) + 32f(a + h/4) + 12f(a + h/2)$$

$$+ 32f(a + 3h/4) + 7f(a + h)\},$$

we have

$$Q_{AS}[A, A + H]f(x) = \sum_{i=1}^{n} Q[a_{i-1}, a_i]f(x).$$

Another interpretation relies on the fact that the five-point Newton-Cotes approximation is the Romberg integration approximation of degree 5 based on two successive subdivisions. In fact, if modification 1 is used to define $\epsilon(r)$, the convergence criterion for each interval is identical with the standard Romberg convergence criterion. That is, the final result for each interval is the same as that which would have been attained if we had set up the Romberg table

$$T_0^{(0)}$$
$$T_1^{(0)}$$
$$T_0^{(1)} \qquad\qquad T_2^{(0)}$$
$$T_1^{(1)}$$
$$T_0^{(2)}$$

for each interval and allowed convergence if $| T_2^{(0)} - T_1^{(1)} | < \epsilon h$. But in this case, if this criterion were not satisfied, instead of calculating the higher degree approximation $T_3^{(0)}$, as would be the case in Romberg integration, we bisect the interval.

It should be noted that the routine having this modification returns a result considerably more accurate than the required $\epsilon$, but there is no readily available criterion to determine a priori how much more accurate it is.


5.  *Round-off Error in Function Values*

In Sections 1–4, the effect of round-off error was disregarded. The author has shown elsewhere [4] that in any automatic routine the effect of round-off error in function values can have a significant bearing on the sequence of calculations undertaken by the routine. In applying the convergence criterion in an interval $[a, b]$, the routine calculates the quantity

$$D[a, b]f(x) = \frac{12}{b - a} [R[a, b]f(x) - R^{(2)}[a, b]f(x)], \tag{5.1}$$

which it then compares with

$$E = 180\epsilon/(B - A). \tag{5.2}$$

(This is equivalent to the comparison (1.10) or (3.6); the use of these quantities rather than those in (3.6) makes it possible to avoid excessive division by 2.)

Expression (5.1) may be expanded in the form

$$D[a, b]f(x) = \tfrac{1}{4}\{f_0 - 4f_1 + 6f_2 - 4f_3 + f_4\}, \tag{5.3}$$

the $f_i$ being regularly spaced function evaluations $f(a + (b - a)i/4)$. As the interval

$[a, b]$ becomes smaller, the reduction in the value of $D[a, b]f(x)$ depends on the cancellation of nearly equal quantities. It is well known (for example in the practice of numerical differentiation) that if an expression such as (5.3) is calculated using values of $f(x)$ which have round-off error of magnitude $\epsilon_{r.0}$, and the true value of $D[a, b]f(x)$ is less than $\epsilon_{r.0}$, the calculated value is quite unrealistic and may be several units of $\epsilon_{r.0}$, with either sign. If it happens that $E < \epsilon_{r.0}$, convergence occurs only if the calculated value of $D[a, b]f(x)$ is precisely zero. Thus under these circumstances this routine might fail to converge and continue to subdivide the interval until a chance cancellation of round-off error gave a value of $D[a, b]f(x)$ which was precisely zero. In an example given in [4], a calculation which should have required about 120 function values in fact used 152,997 function values simply because round-off error prevented it from converging. It should be noted that $\epsilon_{r.0}$ should not be confused with the machine accuracy parameter $\epsilon_m$. If $f(x) = 10^{\pm 6} \cos x$ and the cosine function subroutine is coded to machine accuracy, $\epsilon_{r.0} \cong 10^{\pm 6}\epsilon_m$ is a reasonable estimate. In practice, the user may not know the accuracy of the function, and this accuracy may be different in different ranges of the variable $x$.

In principle, it is easy for this routine to determine the round-off level of the integrand. The routine could calculate a finite-difference table and use this to calculate the overall accuracy of the integrand following standard methods (see, for example, Kopal [3]). However, there is no need to do this explicitly in the ASQ routine.

In the Appendix we prove two Theorems.

THEOREM 1. *If* $f^{(4)}(x) = constant$,

$$D[a, (a + b)/2]f(x) = \tfrac{1}{16}D[a, b]f(x).$$

THEOREM 2. *If* $f(x) \in C^{(4)}[a, b]$ *and* $f^{(4)}(x)$ *is of constant sign in this interval, then*

$$| D[a, (a + b)/2]f(x) | \leq | D[a, b]f(x) |.$$

These theorems may be exploited as follows. Suppose that a particular interval does not converge, the value $D_j = | D[a, b]f(x) |$ being greater than $E$. The routine proceeds to bisect the interval and calculates $D_{j+1} = | D[a, (a + b)/2]f(x) |$. In view of Theorem 1, we should expect that $D_{j+1} \cong D_j/16$ on the average. However, we should certainly expect that $D_{j+1} < D_j$. If we find $D_{j+1} \geq D_j$, one of two possibilities has occurred:

(a) there is a zero of $f^{(4)}(x)$ in the interval $[a, b]$; or

(b) the value of $D_{j+1}$ or $D_j$ is seriously affected by round-off error.

In general, the zeros of $f^{(4)}(x)$ are not of frequent occurrence. Thus, instead of constructing a finite-difference table, the routine may simply compare the current value $D_{j+1}$ with the previous value $D_j$ (if the previous interval did not converge). If the current value is not less than the previous value, then we may conclude that either there is a nearby zero of $f^{(4)}(x)$ or the round-off level has been reached. This leads to a modification using which the initial tolerance $\epsilon$ may be altered by the routine, i.e. changed to a different value $\epsilon'$, as and when it seems appropriate. We define $E' = 180\epsilon'/(B - A)$ in analogy to (5.2).

Modification 4 (unadjusted)

(i) If an interval $j + 1$ has not converged, the previous interval $j$ did not con-

verge, and $D_{j+1} \geq D_j$, then raise the current value of tolerance $E'$ to $E' = D_{j+1}$ (subject to adjustments mentioned below).

(ii) If an interval has converged, and $D_{j+1} \neq 0$, then adjust the current value of tolerance $E'$ such that $E' = \max(E, D_{j+1})$ (subject to adjustments mentioned below).

In practice, if this modification is used as it stands the routine readjusts the tolerance with every small statistical fluctuation of the apparent round-off error level. Also, it does not guard against situation (a) above. Thus it is necessary to arrange a series of different paths for the routine to follow. Some of these are similar to those used in iterative schemes for nonlinear integro-differential equations. In general, a sudden large increase in the level $E'$ is inhibited, and after any moderate increase it is pretended that the interval did not converge. This helps to remedy a situation in which $E'$ is increased simply because $f^{(4)}(x)$ changes sign. On the other hand, a marginal decrease in $E'$ is inhibited. Usually only decreases by a factor of four or more are allowed. Taking into account a large number of different circumstances leads to a lengthy code. It should be noted that the routine enters these sections of code only rarely if at all. In a normal run, in which round-off error is insignificant, the additional cost is simply a single comparison after each interval is considered, i.e. "Is $E = E'$?" or "Is $D_{j+1} < D_j$?" according as the previous interval did or did not converge. The answer is invariably "yes" in such a normal run, and the routine proceeds in the normal manner.

If the tolerance $\epsilon$ is raised and lowered during the calculation, the user may be informed of the cumulative effect of this. If there have been $n - 1$ such adjustments at $x = x_1, x_2, \cdots, x_{n-1}$, the interval has been effectively divided into $n$ subintervals $[x_{i-1}, x_i]$, $i = 1, 2, \cdots, n$ (where $A = x_0$ and $B = x_n$), and a different tolerance $\epsilon_i$ ($i = 1, 2, \cdots, n$) used in each. It is clear that for each subinterval the tolerance is $\epsilon_i (x_i - x_{i-1})/(B - A)$. Thus the tolerance corresponding to the entire interval is

$$\epsilon_{\text{eff}} = \frac{1}{B - A} \sum_{i=1}^{n} \epsilon_i (x_i - x_{i-1}).  \tag{5.4}$$

It is a simple matter for the routine to calculate $\epsilon_{\text{eff}}$. The appropriate quantity is simply added to a running sum each time the tolerance is altered, as well as at the conclusion of the calculation.

Finally, we note that if the tolerance has been raised and the routine is operating near the round-off level of the function, there is no point in forming the fifth-order adjustment, because under these circumstances the function does not approximate at this scale to a polynomial.

We may now complete the statement of modification 4.

Modification 4 (continued)

(iii) The routine returns $\epsilon_{\text{eff}}$ given by (5.4), overwriting the input parameter $\epsilon$.

(iv) If modification 3 is included, the fifth-order adjustment should be included only for sections for which $\epsilon' = \epsilon$.

(v) Modification 1 should be applied only for sections for which $\epsilon' = \epsilon$.

In a routine having modification 4, it is possible for the user to specify $\epsilon = 0$. In this case, the routine produces a "best possible result" and returns an estimate of its accuracy.

## 6. General Remark

Four modifications have been suggested here. They are independent of each other and each may be incorporated into any particular adaptive Simpson routine.

Modification 1 is an alternative to the modification of Section 2. The rule for determining $\epsilon(r)$ is changed in order to produce a less accurate result, hopefully still within the required accuracy $\epsilon$, at a cost of fewer function evaluations.

Modification 2, which replaces trisection by bisection, has the same effect. This is accomplished by reducing (on the average) the margin by which the routine satisfies the convergence criterion.

Modification 3 replaces a third-degree result by a fifth-degree result. Thus it improves the actual accuracy (hopefully) without increasing the work required to obtain the result.

Modification 4 guards against round-off error affecting the strategy. In cases where round-off error is significant, the effect is to produce an equally accurate result at a cost of considerably fewer function evaluations.

These modifications are considered by the author to be improvements, and each corresponds to one of the types of change—(i), (ii), and (iii)—mentioned in the introduction.

None of these modifications affects the basic ideas underlying the original adaptive Simpson routine. The present author has written a code which embodies all the modifications, and numerical experiments have shown that over a wide range of circumstances the results have been improvements, as defined in the Introduction, on the original routine. However, this does not imply that in every case these modifications improve the result. Consequently, the author has merely stated what the modifications are, justified them where possible, and explained their general effect. It is up to the user, or the computing institution, to decide which, if any, of these modifications (possibly in the form of user options) are convenient for a particular purpose.

## APPENDIX. The Euler-Maclaurin Expansion

In this appendix we outline the derivation of eq. (3.1), on which the error analysis is based, and the theorems of Section 5, on which the criterion for the determination of round-off error is based.

This derivation is based on the Euler-Maclaurin expansion in its conventional form, which expresses the difference between the $(m + 1)$-point trapezoidal rule approximation

$$T^{(m)}[a, a + h]f(x) = \frac{h}{m} \left\{ \frac{1}{2}f(a) + \frac{1}{2}f(a + h) + \sum_{j=1}^{m-1} f\left(a + \frac{jh}{m}\right) \right\} \quad (A.1)$$

and the integral $I[a, a + h]f(x)$ to which this approximates as follows. If $f(x) \in C^{(n)}[a, a + h]$, then

$$T^{(m)}[a, a + h]f(x) - I[a, a + h]f(x) = \sum_{j=1}^{q-1} \frac{B_j}{j!} \frac{h^j}{m^j} \int_a^{a+h} f^{(j)}(x)\,dx$$

$$+ \frac{h^q}{m^q} \int_a^{a+h} \phi_q(m(x - a)/h)f^{(q)}(x)\,dx, \quad q \le n, \quad (A.2)$$

where $B_j = B_j(1)$ is the $j$th Bernoulli number (0 if $j$ is odd), $B_1$ is set equal to zero, and

$$\phi_j(x) = (B_j - \bar{B}_j(1 - x))/j!. \tag{A.3}$$

Here the function $\bar{B}_j(x)$ is periodic in $x$ with period 1 and coincides with the $j$th Bernoulli polynomial $B_j(x)$ in the interval $0 < x < 1$. These functions and numbers are given in Abramowitz and Stegun [1, pp. 803 et seq].

The corresponding formula for Simpson's $(2m + 1)$-point rule is obtained by noting that

$$R^{(m)}[a, a + h]f(x) = \tfrac{4}{3}T^{(2m)}[a, a + h]f(x) - \tfrac{1}{3}T^{(m)}[a, a + h]f(x). \tag{A.4}$$

Both terms on the right-hand side may be replaced using (A.2) to give

$$R^{(m)}[a, a + h]f(x) - I[a, a + h]f(x) = \sum_{r=1}^{q-1} c_r \frac{h^r}{m^r} \int_a^{a+h} f^{(r)}(x)\, dx$$
$$+ \frac{h^q}{m^q} \int_a^{a+h} g_q(m(x - a)/h)f^{(q)}(x)\, dx, \qquad q \le n, \tag{A.5}$$

where

$$c_r = \left(\frac{4}{3}2^{-r} - \frac{1}{3}\right)\frac{B_r}{r!} \tag{A.6}$$

and

$$g_q(x) = \tfrac{4}{3}2^{-q}\phi_q(2x) - \tfrac{1}{3}\phi_q(x). \tag{A.7}$$

We note that $c_r = 0$ ($r$ odd) and $c_2 = 0$. Equation (3.1) is obtained by setting $q = 5$ in eq. (A.5).

Theorems 1 and 2 of Section 5 are concerned with a difference calculated in the routine, namely,

$$D[a, a + h]f(x) = \frac{12}{h}[R[a, a + h]f(x) - R^{(2)}[a, a + h]f(x)]. \tag{A.8}$$

If $f^{(4)}(x) = K$ is a constant, eq. (3.2) gives

$$D[a, a + h]f(x) = \tfrac{15}{16}\cdot 12c_4Kh^4,$$

and Theorem 1 follows directly.

We use (A.5) with $q = 4$ and $m = 1$, and $m = 2$, to derive Theorem 2. Thus

$$D_1 \equiv D[a, a + h]f(x) = 12h^3 \int_a^{a+h} G_4((x - a)/h)f^{(4)}(x)\, dx, \tag{A.9}$$

$$D_1 - D_2 - D_3 \equiv D[a, a + h]f(x) - D[a, a + h/2]f(x)$$
$$-D[a + h/2, a + h]f(x) = 12h^3 \int_a^{a+h} F_4((x - a)/h)f^{(4)}(x)\, dx, \tag{A.10}$$

where

$$G_4(x) = g_4(x) - 2^{-4}g_4(2x), \tag{A.11}$$

$$F_4(x) = G_4(x) - 2^{-3}G_4(2x). \tag{A.12}$$

The proof of Theorem 2 depends on showing successively that $g_4(x)$, $G_4(x)$, and $F_4(x)$ are nonnegative definite. These are all spline polynomials of degree 4 and may be evaluated explicitly using eqs. (A.3), (A.7), (A.11), and (A.12). A tedious but straightforward calculation yields:

$$4!\phi_4(x) = -x^2(x-1)^2, \qquad 0 < x < 1,$$

$$4!g_4(x) = -x^3(3x-2)/3, \qquad 0 < x < \tfrac{1}{2},$$

$$4!G_4(x) = x^3/3, \qquad 0 < x < \tfrac{1}{4},$$

$$= [1 - 12x(2x-1)^2]/48, \qquad \tfrac{1}{4} < x < \tfrac{1}{2},$$

$$4!F_4(x) = 0, \qquad 0 < x < \tfrac{1}{8},$$

$$= (8x-1)^3/384, \qquad \tfrac{1}{8} < x < \tfrac{1}{4},$$

$$\geq 4!F_4(\tfrac{1}{4}), \qquad \tfrac{1}{4} < x < \tfrac{1}{2}.$$

These functions are all of period 1 and are symmetric about $x = \tfrac{1}{2}$. The latter three functions are nowhere negative.

A condition of the theorem is that $f^{(4)}(x)$ is of the same sign (or zero) throughout the interval $[a, a + h]$. For definiteness, we take this sign to be positive. Then it follows from eqs. (A.9) and (A.10) and the definite signs of $F_4(x)$ and $G_4(x)$ that

$$D_1 \geq 0, \qquad D_2 \geq 0, \qquad D_3 \geq 0, \qquad D_1 - D_2 - D_3 \geq 0.$$

These inequalities can be satisfied only if $D_1 \geq D_2$. This establishes Theorem 2 in the case $f^{(4)}(x)$ is nonnegative. A trivial modification of the argument establishes $|D_1| \geq |D_2|$ in the case $f^{(4)}(x)$ is nonpositive.

REFERENCES

1. ABRAMOWITZ, M., AND STEGUN, I. A. *Handbook of Mathematical Functions*. Dover, New York, 1965.
2. DAVIS, P., AND RABINOWITZ, P. *Numerical Integration*. Blaisdell, London, 1967.
3. KOPAL, Z. *Numerical Analysis*. Chapman and Hall, London, 1955.
4. LYNESS, J. N. The effect of inadequate convergence criteria in automatic routines. (To appear in *Comput. J.*)
5. ——, AND NINHAM, B. W. Asymptotic expansions and numerical quadrature. *Math. Comput. 21* (1967), 162–178.
6. MCKEEMAN, W. M. Algorithm 145, adaptive numerical integration by Simpson's rule. *Comm. ACM 5*, 12 (Dec. 1962), 604.
7. ——. Certification of algorithm 145; adaptive numerical integration by Simpson's rule. *Comm. ACM 6*, 4 (Apr. 1963), 167–168.
8. —— AND TESLER, L. Algorithm 182, nonrecursive adaptive integration. *Comm. ACM 6*, 6 (June 1963), 315.
9. ——. Algorithm 198, adaptive integration and multiple integration. *Comm. ACM 6*, 8 (Aug. 1963), 443.