

On the Round-Off Error in the Method of Successive Over-Relaxation

By M. Stuart Lynn

Abstract. The asymptotic behavior of the round-off error, which accumulates when the well-known iterative method of (point) successive over-relaxation is used to solve a large-scale system of linear equations, is examined by means of a statistical model. The local round-off errors are treated as independent random variables and expressions for the mean and variance of the accumulated round-off error are obtained, as the number of iterations tends to infinity.

1. Introduction. Consider the system of linear equations

$$(1.1) \quad \mathbf{C}\mathbf{x} = \mathbf{b},$$

where \mathbf{C} is an $n \times n$ real matrix and \mathbf{b} an $n \times 1$ vector. In the numerical solution of (1.1), round-off errors accumulate. Wilkinson [14, 15] and Turing [11] have considered the effect of this where direct methods of solution are involved. As far as iterative methods are concerned, bounds for the round-off error occurring in a general iterative procedure have been obtained by Urabe [12] and Descloux [2].

In order to attempt to obviate the usual criticism that such bounds may be unrealistic, we shall use the technique, familiar in other fields of numerical analysis, e.g. Henrici [7], of treating the local round-off errors (see below) as independent random variables, and then obtaining expressions for the asymptotic behaviour of the mean and variance of the accumulated round-off error. In connection with iterative procedures, the statistical model employed was first used by Abramov [1] in studying the round-off error generated by the Jacobi procedure for solving (1.1). More recently, Golub [5] has used similar techniques to analyze the Richardson second-order and Chebyshev semi-iterative methods. We shall consider another method, namely the method of (point) successive over-relaxation: Young [16], Varga [13].

The statistical model which is employed is certainly open to many severe criticisms, e.g. Forsythe [3], more so, in fact, than in other fields of numerical analysis where perhaps it can be used with more confidence. Principally, the fundamental assumption that the local round-off errors are either independent or random is certainly questionable, particularly after a large number of iterations, and the procedure reaches a 'state of numerical convergence': Sibagaki [10], Urabe [12]. However, the somewhat curious nature of the results obtained from this kind of analysis (see particularly Theorem 2) are perhaps not without interest. We would again refer the reader to Golub [4, 5] for a more detailed appraisal of the use of a statistical model in connection with iterative techniques, compared with the use of bounds.

Received February 25, 1963. Revised May 20, 1963. This work was carried out mainly under the sponsorship of the Office of Naval Research, and partly within the research program of the National Physical Laboratory. The paper is published by permission of the Director of the Laboratory.

Having obtained expressions for the mean and variance of the accumulated round-off error, it would not be strictly correct to apply the central limit theorem and deduce a normal distribution, since it can readily be shown that the accumulated round-off error remains uniformly bounded as the number of iterations tends to infinity. However, using an analysis precisely analogous to that used by Golub [5] we can nevertheless obtain probabilistic bounds for the round-off error. The contents of this are expressed in Theorem 3.

The method of point successive over-relaxation or the Young-Frankel method may, as usual, be defined as follows. Let \mathbf{C} be decomposed by

$$(1.2) \quad \mathbf{C} = \mathbf{S} - \mathbf{L} - \mathbf{U} = \mathbf{S} - \mathbf{B}$$

where \mathbf{S} is a diagonal matrix and \mathbf{L} , \mathbf{U} are strictly lower and upper triangular matrices, respectively. For arbitrary $\mathbf{x}^{(0)}$, the sequence $\{\mathbf{x}^{(k)}\}$ of vectors is defined by

$$(1.3) \quad \mathbf{S}\mathbf{x}^{(k+1)} = \omega(\mathbf{L}\mathbf{x}^{(k+1)} + \mathbf{U}\mathbf{x}^{(k)} + \mathbf{b}) + (1 - \omega)\mathbf{S}\mathbf{x}^{(k)} \quad k = 0, 1, 2, \dots,$$

where ω is a *relaxation-parameter* which necessarily lies in the range

$$(1.4) \quad 0 < \omega < 2.$$

We shall henceforth assume that \mathbf{C} is symmetric positive-definite, so that

$$(1.5) \quad \mathbf{C} = \mathbf{C}^*,$$

and

$$(1.6) \quad \mathbf{L}^* = \mathbf{U}, \quad \mathbf{B} = \mathbf{B}^*,$$

where, in general, \mathbf{A}^* denotes the (conjugate) transpose of a matrix \mathbf{A} (since we shall always be working with real matrices, the conjugation operation is not involved). Condition (1.4) is then both necessary and sufficient to ensure convergence of (1.3) to \mathbf{x} as $k \rightarrow \infty$: Ostrowski [9]. Without loss of generality, we may also assume that

$$(1.7) \quad \mathbf{S} = \mathbf{I}$$

the identity matrix, since we may write (1.1) in the form,

$$(1.8) \quad \mathbf{S}^{-1/2}\mathbf{A}\mathbf{S}^{-1/2}\mathbf{S}^{1/2}\mathbf{x} = \mathbf{S}^{-1/2}\mathbf{b},$$

or

$$(1.9) \quad \mathbf{A}_1\mathbf{y} = \mathbf{b}_1,$$

where

$$(1.10) \quad \mathbf{A}_1 = \mathbf{S}^{-1/2}\mathbf{A}\mathbf{S}^{-1/2}, \quad \mathbf{y} = \mathbf{S}^{1/2}\mathbf{x}, \quad \mathbf{b}_1 = \mathbf{S}^{-1/2}\mathbf{b},$$

and \mathbf{A}_1 is symmetric positive-definite with unity diagonal entries.

2. The Model. Now suppose that, instead of solving (1.3) exactly, we actually compute vectors $\tilde{\mathbf{x}}^{(k)}$, where (replacing \mathbf{S} by \mathbf{I})

$$(2.1) \quad \tilde{\mathbf{x}}^{(k+1)} = \omega(\mathbf{L}\tilde{\mathbf{x}}^{(k+1)} + \mathbf{U}\tilde{\mathbf{x}}^{(k)} + \mathbf{b}) + (1 - \omega)\tilde{\mathbf{x}}^{(k)} + \mathbf{e}^{(k+1)}, \quad k = 0, 1, 2, \dots$$

Here, $\{\mathbf{e}^{(k)}\} (k = 0, 1, 2, \dots)$ are the *local round-off errors* in the sense that they are the errors committed at each iteration. Let

$$(2.2) \quad \mathbf{r}^{(k)} = \bar{\mathbf{x}}^{(k)} - \mathbf{x}^{(k)} \quad (k = 0, 1, 2, \dots)$$

be the *accumulated round-off error vector* at the k^{th} iteration. Then from (1.3) and (2.1),

$$(2.3) \quad \mathbf{r}^{(k+1)} = \omega(\mathbf{L}\mathbf{r}^{(k+1)} + \mathbf{U}\mathbf{r}^{(k)}) + (1 - \omega)\mathbf{r}^{(k)} + \mathbf{e}^{(k+1)}, \quad k = 0, 1, 2, \dots$$

where

$$(2.4) \quad \mathbf{r}^{(0)} = \mathbf{e}^{(0)} = \mathbf{0}$$

We treat the $\{\mathbf{e}^{(k)}\}$ as independent random variables (see the remarks in §1) with expected value

$$E[\mathbf{e}^{(k)}] = \boldsymbol{\varepsilon}_k,$$

where we assume that

$$(2.5) \quad \boldsymbol{\varepsilon}_k \equiv \boldsymbol{\varepsilon} \quad (k > 0)$$

where $\boldsymbol{\varepsilon}$ is constant; and with co-variance

$$(2.6) \quad D\mathbf{e}^{(k)} = \sigma^2\mathbf{R} \quad (k > 0); \quad D\mathbf{e}^{(0)} = \mathbf{0}$$

where \mathbf{R} is some positive definite symmetric matrix, again assumed constant, which we further necessarily assume commutes with \mathbf{L} and \mathbf{U} , and hence with \mathbf{A} and \mathbf{B} . By definition, we have that

$$(2.7) \quad \begin{aligned} D\mathbf{e}^{(k)} &= E[(\mathbf{e}^{(k)} - \boldsymbol{\varepsilon})(\mathbf{e}^{(k)} - \boldsymbol{\varepsilon})^*] \\ &= E[\mathbf{e}^{(k)}\mathbf{e}^{(k)*}] - \boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^*, \end{aligned}$$

and so

$$(2.8) \quad E[\mathbf{e}^{(k)}\mathbf{e}^{(k)*}] = \sigma^2\mathbf{R} + \boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^*$$

Since we also assume that the $\{\mathbf{e}^{(k)}\}$ are independent, we also have that

$$(2.9) \quad \begin{aligned} E[\mathbf{e}^{(k)}\mathbf{e}^{(j)*}] &= E[\mathbf{e}^{(k)}]E[\mathbf{e}^{(j)*}] \\ &= \boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^* \quad (k \neq j) \end{aligned}$$

Let

$$(2.10) \quad \mathbf{r}_k = E[\mathbf{r}^{(k)}]$$

denote the expected-value of the accumulated round-off error at the k^{th} iteration, and let

$$\mathbf{V}_k = D\mathbf{r}^{(k)}$$

denote the co-variance matrix of the same. We are interested in the behavior of \mathbf{r}_k and \mathbf{V}_k as $k \rightarrow \infty$.

Now from (2.3)

$$\mathbf{r}^{(k+1)} = \mathbf{T}\mathbf{r}^{(k)} + (\mathbf{I} - \omega\mathbf{L})^{-1}\mathbf{e}^{(k+1)},$$

where

$$\mathbf{T} = (\mathbf{I} - \omega\mathbf{L})^{-1}(\omega\mathbf{U} + \overline{\mathbf{I} - \omega\mathbf{I}}).$$

Thus

$$(2.11) \quad \mathbf{r}^{(k+1)} = \sum_{j=1}^{k+1} \mathbf{T}^{k+1-j} (\mathbf{I} - \omega\mathbf{L})^{-1} \mathbf{e}^{(j)},$$

and so, since E is a linear operator,

$$\begin{aligned} \mathbf{r}_{k+1} &= \left\{ \sum_{j=1}^{k+1} \mathbf{T}^{k+1-j} \right\} (\mathbf{I} - \omega\mathbf{L})^{-1} \boldsymbol{\varepsilon} \\ &= (\mathbf{I} - \mathbf{T}^{k+1}) (\mathbf{I} - \mathbf{T})^{-1} (\mathbf{I} - \omega\mathbf{L})^{-1} \boldsymbol{\varepsilon} \\ &= (\mathbf{I} - \mathbf{T}^{k+1}) \omega^{-1} \mathbf{C}^{-1} \boldsymbol{\varepsilon}. \end{aligned}$$

Hence we have:

THEOREM 1. *Under the foregoing statistical assumptions, and if ω satisfies (1.4), then:*

$$\mathbf{r}_k = (\mathbf{I} - \mathbf{T}^k) \omega^{-1} \mathbf{C}^{-1} \boldsymbol{\varepsilon},$$

and

$$(2.12) \quad \mathbf{r} = \frac{1}{\omega} \mathbf{C}^{-1} \boldsymbol{\varepsilon},$$

where

$$(2.13) \quad \mathbf{r} = \lim_{k \rightarrow \infty} \mathbf{r}_k.$$

It should be noted that Theorem 1 does not rely upon the assumed independence of the $\{\mathbf{e}^{(k)}\}$.

3. The Co-Variance Matrix. We now consider the behavior of \mathbf{V}_k as $k \rightarrow \infty$. From (2.8), (2.10) and the commutativity assumptions, we have immediately that†

$$(3.1) \quad \mathbf{V}_k = \sigma^2 \mathbf{R} [(\mathbf{I} - \omega\mathbf{U})(\mathbf{I} - \omega\mathbf{L})]^{-1} \sum_{j=1}^k \mathbf{T}^{k-j} \mathbf{T}^{*k-j}.$$

Since \mathbf{T} is neither symmetric nor normal, we lack any method for finding an explicit expression for the right-hand side of (3.1) (bounds do not interest us, since they defeat the object of a statistical investigation). This is unlike the case for Jacobi iteration as studied by Abramov [1], and it is in this context that the present problem possesses a separate interest. This is similarly unlike the situation for the Richardson second-order method as studied by Golub [5].

For a particular class of matrices, however, we may adopt a different approach; this is the class of matrices which have Property A and are σ_1 -ordered: Young [16]. For our purposes we may define these as matrices for which \mathbf{B} in (1.2) has

† The author is indebted to the referee for correcting an error in the author's original analysis, and for the more fruitful formulation and proof of Theorem 1.

the form:

$$(3.2) \quad \mathbf{B} = m \left\{ \begin{array}{c|c} \widetilde{0} & F \\ \hline F^* & 0 \end{array} \right\}^m$$

for some integer m , $0 < m < n$, so that F is an $m \times (n - m)$ sub-matrix. Thus we have:

Definition 2.1. If \mathbf{B} in (1.2) has the form (3.2), we shall say that \mathbf{C} in (1.2) has Property A and is σ_1 -ordered.

In the terminology of Varga [13] \mathbf{B} is weakly cyclic of index 2 and is in its normal form.

It is also worth remarking, perhaps, that if $\mathbf{C} = \mathbf{C}^*$, has Property A and is σ_1 -ordered, then the assumption that \mathbf{R} commutes with both \mathbf{L} and \mathbf{U} is satisfied if (i) \mathbf{R} has the form:

$$(3.3) \quad \mathbf{R} = m \left\{ \begin{array}{c|c} \widetilde{R_1} & 0 \\ \hline 0 & R_2 \end{array} \right\}^m$$

where $R_1 = R_1^*$ is an $m \times m$ sub-matrix, and $R_2 = R_2^*$ is an $(n - m) \times (n - m)$ sub-matrix, both of which are, of course, positive-definite; and (ii) $FR_1 = R_2F$.

Our principal result is:

THEOREM 2. *If $\mathbf{C} = \mathbf{C}^*$ is symmetric positive-definite, has Property A and is σ_1 -ordered, if $0 < \omega < 2$, and if (2.5), (2.6) and (2.9) are satisfied, then*

$$(3.4) \quad \mathbf{V} = \sigma^2 / [\omega(2 - \omega)] \cdot \mathbf{R}\mathbf{C}^{-1},$$

where

$$(3.5) \quad \mathbf{V} = \lim_{k \rightarrow \infty} \mathbf{V}_k,$$

and

$$(3.6) \quad \mathbf{V}_k = D\mathbf{r}^{(k)},$$

the co-variance of $\mathbf{r}^{(k)}$.

Proof. The proof proceeds by a sequence of lemmas. As an immediate consequence of (3.2), we find that:

LEMMA 1. *If \mathbf{C} has Property A and is σ_1 -ordered, then*

$$(3.7) \quad \mathbf{L}^2 = \mathbf{U}^2 = \mathbf{0},$$

and hence \mathbf{B}^2 commutes with both \mathbf{L} and \mathbf{U} .

LEMMA 2. *If $\mathbf{C} = \mathbf{C}^*$ has Property A and is σ_1 -ordered, then the method defined by (1.3) converges if and only if the roots of*

$$(3.8) \quad \theta^2 - \lambda_i \theta - \alpha = 0 \quad (i = 1, \dots, n)$$

are less than one in modulus, where λ_i are the eigenvalues of

$$(3.9) \quad \mathbf{A} = 2(1 - \omega)\mathbf{I} + \omega^2\mathbf{B}^2,$$

and

$$(3.10) \quad \alpha = (1 - \omega)^2.$$

Proof. This is merely a re-formulation of the familiar relationship of Young [16], equation (2.4). We note that, for future reference,

$$(3.11) \quad 0 \leq \alpha < 1$$

LEMMA 3. $\{\mathbf{r}^{(k)}\} (k = 0, 1, 2, \dots)$ satisfy the relationships:

$$(3.12) \quad \mathbf{r}^{(k+2)} = \mathbf{A}\mathbf{r}^{(k+1)} - \alpha\mathbf{r}^{(k)} + \mathbf{f}^{(k+2)} \quad (k = 0, 1, 2, \dots)$$

$$(3.13) \quad \mathbf{r}^{(0)} = \mathbf{f}^{(0)} = \mathbf{0}; \quad \mathbf{r}^{(1)} = \mathbf{f}^{(1)} = (\mathbf{I} + \omega\mathbf{L})\mathbf{e}^{(1)},$$

where

$$(3.14) \quad \mathbf{f}^{(k+2)} = (\mathbf{I} + \omega\mathbf{L})\mathbf{e}^{(k+2)} + (\omega\mathbf{U} + \overline{\mathbf{I} - \omega\mathbf{I}})\mathbf{e}^{(k+1)} \quad (k = 0, 1, 2, \dots)$$

Proof. Since $\mathbf{L}^2 = \mathbf{0}$, one can immediately verify that

$$(\mathbf{I} - \omega\mathbf{L})^{-1} = (\mathbf{I} + \omega\mathbf{L}).$$

Hence, from (2.3), for $k = 0, 1, 2, 3, \dots$,

$$\begin{aligned} \mathbf{r}^{(k+2)} &= (\mathbf{I} - \omega\mathbf{L})^{-1}[(\omega\mathbf{U} + \overline{\mathbf{I} - \omega\mathbf{I}})\mathbf{r}^{(k+1)} + \mathbf{e}^{(k+2)}] \\ &= (\mathbf{I} + \omega\mathbf{L})[(\omega\mathbf{U} + \overline{\mathbf{I} - \omega\mathbf{I}})\mathbf{r}^{(k+1)} + \mathbf{e}^{(k+2)}] \\ &= (\overline{\mathbf{I} - \omega\mathbf{I}} + \omega^2\mathbf{B}^2)\mathbf{r}^{(k+1)} + \omega\mathbf{U}(\mathbf{I} - \omega\mathbf{L})\mathbf{r}^{(k+1)} + \omega(1 - \omega)\mathbf{L}\mathbf{r}^{(k+1)} \\ &\quad + (\mathbf{I} + \omega\mathbf{L})\mathbf{e}^{(k+2)} \\ &= (\overline{\mathbf{I} - \omega\mathbf{I}} + \omega^2\mathbf{B}^2)\mathbf{r}^{(k+1)} + \omega\mathbf{U}[\omega\mathbf{U}\mathbf{r}^{(k)} + (1 - \omega)\mathbf{r}^{(k)} + \mathbf{e}^{(k+1)}] \\ &\quad + \omega(1 - \omega)\mathbf{L}\mathbf{r}^{(k+1)} + (\mathbf{I} + \omega\mathbf{L})\mathbf{e}^{(k+2)} \\ &= (\overline{\mathbf{I} - \omega\mathbf{I}} + \omega^2\mathbf{B}^2)\mathbf{r}^{(k+1)} + \omega(1 - \omega)[\mathbf{L}\mathbf{r}^{(k+1)} + \mathbf{U}\mathbf{r}^{(k)}] \\ &\quad + (\mathbf{I} + \omega\mathbf{L})\mathbf{e}^{(k+2)} + \omega\mathbf{U}\mathbf{e}^{(k+1)} \\ &= (2\overline{\mathbf{I} - \omega\mathbf{I}} + \omega^2\mathbf{B}^2)\mathbf{r}^{(k+1)} - (1 - \omega)^2\mathbf{r}^{(k)} + \mathbf{f}^{(k+2)} \\ &= \mathbf{A}\mathbf{r}^{(k+1)} - \alpha\mathbf{r}^{(k)} + \mathbf{f}^{(k+2)} \end{aligned}$$

as required. The proof of (3.13) is entirely similar.

Now from (3.12) and (3.13), it follows that for certain polynomials $p_{kj}(\mathbf{A})$ in \mathbf{A} ,

$$(3.15) \quad \mathbf{r}^{(k)} = \sum_{j=1}^k p_{kj}(\mathbf{A})\mathbf{f}^{(j)} \quad (k = 1, 2, 3, \dots)$$

Hence

$$\begin{aligned} \mathbf{r}_k &= E[\mathbf{r}^{(k)}] = \sum_{j=1}^k p_{kj}(\mathbf{A})E[\mathbf{f}^{(j)}] \quad (k = 1, 2, 3, \dots) \\ \mathbf{r}_0 &= \mathbf{0}; \quad \mathbf{r}_1 = E[\mathbf{f}^{(1)}] = (\mathbf{I} + \omega\mathbf{L})\mathbf{e} \end{aligned}$$

Thus, since $\mathbf{A} = \mathbf{A}^*$,

$$(3.17) \quad \begin{aligned} \mathbf{V}_k &= E[\mathbf{r}^{(k)}\mathbf{r}^{(k)*}] - \mathbf{r}_k\mathbf{r}_k^* \\ &= \sum_{i,j=1}^k p_{ki}(\mathbf{A}) \mathbf{E}_{ij} p_{kj}(\mathbf{A}) \end{aligned}$$

where

$$(3.18) \quad \mathbf{E}_{ij} = E[\mathbf{f}^{(i)}\mathbf{f}^{(j)*}] - E[\mathbf{f}^{(i)}]E[\mathbf{f}^{(j)*}]$$

We now wish to compute the $\mathbf{E}_{ij}(i, j = 1, 2, \dots)$. Various cases arise. From (2.5), (2.8), (2.9) and the fact that \mathbf{R} commutes with \mathbf{B} , \mathbf{L} and \mathbf{U} , we have, omitting details:

$$(3.19) \quad \begin{aligned} \text{(i) } i = j \neq 1: & \quad \mathbf{E}_{i,i} = \sigma^2\mathbf{R}[2(1 - \omega)\mathbf{I} + \omega^2(\mathbf{I} + \mathbf{B} + \mathbf{B}^2)] \\ \text{(ii) } i = j = 1: & \quad \mathbf{E}_{1,1} = \sigma^2\mathbf{R}[\mathbf{I} + \omega\mathbf{B} + \omega^2\mathbf{L}\mathbf{U}] \\ \text{(iii) } 2 \leq j = i + 1 \leq k: & \quad \mathbf{E}_{i,i+1} = \sigma^2\mathbf{R}[\omega^2\mathbf{L} - (1 - \omega)\mathbf{I}] \\ \text{(iv) } 2 \leq i = j + 1 \leq k: & \quad \mathbf{E}_{i+1,i} = \sigma^2\mathbf{R}[\omega^2\mathbf{U} - (1 - \omega)\mathbf{I}] \\ \text{(v) } |i - j| > 2: & \quad \mathbf{E}_{i,j} = \mathbf{0} \end{aligned}$$

Let $\{\mu_i\}(i = 1, \dots, n)$ be the real eigenvalues of \mathbf{B} , so that

$$\lambda_i = 2(1 - \omega) + \omega^2\mu_i^2 \quad (i = 1, \dots, n)$$

where $\lambda_i(i = 1, \dots, n)$ are defined in the statement of Lemma 2. Now, if $p(\mathbf{B})$ is any polynomial in \mathbf{B} , then the eigenvalues of $p(\mathbf{B})$ are $p(\mu_i)(i = 1, \dots, n)$ and the corresponding eigenvectors are precisely the n linearly independent, orthonormal (say) eigenvectors $\mathbf{x}_i(i = 1, \dots, n)$ of \mathbf{B} . Then from (3.17) and (3.19), and using the spectral decomposition of \mathbf{B} and \mathbf{A} ,

$$(3.20) \quad \begin{aligned} \mathbf{V}_k &= \sigma^2\mathbf{R} \left\{ [2(1 - \omega)\mathbf{I} + \omega^2(\mathbf{I} + \mathbf{B} + \mathbf{B}^2)] \cdot \sum_{i=2}^k p_{k,i}^2(\mathbf{A}) \right. \\ &\quad + [\mathbf{I} + \omega\mathbf{B} + \omega^2\mathbf{L}\mathbf{U}] p_{k,1}^2(\mathbf{A}) \\ &\quad \left. + [\omega^2\mathbf{B} - 2(1 - \omega)\mathbf{I}] \cdot \sum_{i=1}^{k-1} p_{k,i}(\mathbf{A}) p_{k,i-1}(\mathbf{A}) \right\} \\ &= \sigma^2\mathbf{R} \cdot \sum_{j=1}^n \left\{ [2(1 - \omega) + \omega^2(1 + \mu_j + \mu_j^2)] \cdot \sum_{i=2}^k p_{k,i}^2(\lambda_j) \right. \\ &\quad + [\mathbf{I} + \omega\mathbf{B} + \omega^2\mathbf{L}\mathbf{U}] p_{k,1}^2(\lambda_j) \\ &\quad \left. + [\omega^2\mu_j - 2(1 - \omega)] \cdot \sum_{i=1}^{k-1} p_{k,i}(\lambda_j) p_{k,i-1}(\lambda_j) \right\} \mathbf{x}_j \mathbf{x}_j^* \end{aligned}$$

where $p_{k,j}(\varphi)$, for any scalar, φ , denotes that polynomial in φ which has the same coefficients as $p_{k,j}(\mathbf{A})$.

In order to examine the behaviour of (3.20) as $k \rightarrow \infty$, we must consider the $p_{k,j}$ in more detail. Substituting (3.15) into (3.12), (3.13) and (3.14), and 'comparing coefficients' of the $\{\mathbf{f}^{(k)}\}$ (which is permissible since the resulting expression is to be true for *all* $\{\mathbf{e}^{(k)}\}$ and hence *all* $\{\mathbf{f}^{(k)}\}$), we have the recurrence relationships:

$$\begin{aligned}
(3.21) \quad & p_{k,k}(\mathbf{A}) = \mathbf{I} \\
& p_{k,k-1}(\mathbf{A}) = \mathbf{A} \\
& p_{k,i}(\mathbf{A}) - \mathbf{A}p_{k-1,i}(\mathbf{A}) + \alpha p_{k-2,i}(\mathbf{A}) = \mathbf{0} \quad (i = 1, 2, \dots, k-2)
\end{aligned}$$

for $k = 1, 2, \dots$. Similarly, $p_{k,j}(\varphi)$ satisfies the recurrence relationships

$$\begin{aligned}
(3.22) \quad & p_{k,k}(\varphi) = 1 \\
& p_{k,k-1}(\varphi) = \varphi \\
& p_{k,i}(\varphi) - \varphi p_{k-1,i}(\varphi) + \alpha p_{k-2,i}(\varphi) = 0 \quad (i = 1, 2, \dots, k-2)
\end{aligned}$$

for $k = 1, 2, \dots$. Now consider the difference equation for $p_i(\varphi)$

$$\begin{aligned}
(3.23) \quad & p_{i+2}(\varphi) - \varphi p_{i+1}(\varphi) + \alpha p_i(\varphi) = 0 \quad (i = 0, 1, 2, \dots) \\
& p_0(\varphi) = 1; \quad p_1(\varphi) = \varphi
\end{aligned}$$

Then from the familiar theory of difference equations with constant coefficients, we know that

$$(3.24) \quad p_i(\varphi) = \begin{cases} (t_1^{i+1} - t_2^{i+1}) / (t_1 - t_2), & \text{if } t_1 \neq t_2 \\ (i+1)t_1^i, & \text{if } t_1 = t_2 \end{cases}$$

where t_1, t_2 are the roots of

$$(3.25) \quad t^2 - \varphi t + \alpha = 0$$

From (3.23) we see that $p_i(\varphi)$ is a polynomial in φ . It is easy to verify from (3.22), and noting that both the $p_i(\varphi)$ and the $p_{k,j}(\varphi)$ are uniquely defined, that

$$(3.26) \quad p_{k,j}(\varphi) = p_{k-j}(\varphi) \quad (j = 1, \dots, k)$$

Hence

$$(3.27) \quad p_{k,j}(\mathbf{A}) = p_{k-j}(\mathbf{A})$$

where $p_i(\mathbf{A})$ has the obvious meaning. It follows that

$$(3.28) \quad \sum_{i=2}^k p_{k,i}^2(\varphi) = \sum_{i=2}^k p_{k-i}^2(\varphi) = \sum_{i=0}^{k-2} p_i^2(\varphi)$$

$$(3.29) \quad \sum_{i=2}^{k-1} p_{k,i}(\varphi)p_{k,i+1}(\varphi) = \sum_{i=2}^{k-1} p_{k-i}(\varphi)p_{k-i-1}(\varphi) = \sum_{i=1}^{k-2} p_i(\varphi)p_{i-1}(\varphi)$$

Furthermore, from (3.24), if $|t_s| < 1$ ($s = 1, 2$), then $p_i(\varphi) \rightarrow 0$ as $i \rightarrow \infty$, and hence from Lemma 2,

$$(3.30) \quad p_{k,1}(\lambda_j) = p_{k-1}(\lambda_j) \rightarrow 0 \quad \text{as } k \rightarrow \infty \quad j = 1, \dots, n.$$

We now need:

LEMMA 4. Suppose, in (3.25), that $|t_s| < 1$ ($s = 1, 2$), and let

$$(3.31) \quad M_k(\varphi) = \sum_{i=0}^k p_i^2(\varphi); \quad N_k(\varphi) = \sum_{i=1}^k p_i(\varphi)p_{i-1}(\varphi).$$

Then

$$(3.32) \quad M(\varphi) = \lim_{k \rightarrow \infty} M_k(\varphi), \quad \text{and} \quad N(\varphi) = \lim_{k \rightarrow \infty} N_k(\varphi)$$

both exist, and in fact,

$$(3.33) \quad \begin{aligned} M(\varphi) &= [(1 + \alpha)/(1 - \alpha)]/[(1 + \alpha)^2 - \varphi^2] \\ N(\varphi) &= [\varphi/(1 - \alpha)]/[(1 + \alpha)^2 - \varphi^2] \end{aligned}$$

Proof. This is a simple computation using (3.24), the fact $t_1 + t_2 = \varphi$, $t_1 t_2 = \alpha$ and that $|t_s| < 1$ ($s = 1, 2$).

Finally, then, applying Lemma 2 in Lemma 4, and using this in (3.28), (3.29) and (3.30), we have, on letting $k \rightarrow \infty$ in (3.20):

$$(3.34) \quad \begin{aligned} \mathbf{V} &= \lim_{k \rightarrow \infty} \mathbf{V}_k \\ &= \sigma^2 \mathbf{R} \sum_{j=1}^n \{ [2(1 - \omega) + \omega^2(1 + \mu_j + \mu_j^2)] \cdot [1 + (1 - \omega)^2] / [1 - (1 - \omega)^2] \\ &\quad + \lambda_j [\omega^2 \mu_j - 2(1 - \omega)] / [1 - (1 - \omega)^2] \} \frac{\mathbf{x}_j \mathbf{x}_j^*}{[1 + (1 - \omega)^2] - \lambda_j^2} \\ &= \sigma^2 \mathbf{R} \sum_{j=1}^n Q_j(\omega) \mathbf{x}_j \mathbf{x}_j^*, \end{aligned}$$

where, since $\lambda_j = 2(1 - \omega) + \omega^2 \mu_j^2$, and simplifying:

$$(3.35) \quad \begin{aligned} Q_j(\omega) &= \frac{[\omega^4 \mu_j^3 + \omega^4 \mu_j^2 + \omega^2(2 - \omega)^2 \mu_j + \omega^2(2 - \omega)^2]}{\omega^3[(2 - \omega)^2 + \omega^2 \mu_j^2](1 - \mu_j^2)} \\ &= 1/[\omega(2 - \omega)(1 - \mu_j)]. \end{aligned}$$

Thus

$$\mathbf{V} = \sigma^2 \mathbf{R} / [\omega(2 - \omega)] \cdot \sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^* / (1 - \mu_j) = \sigma^2 / [\omega(2 - \omega)] \cdot \mathbf{R} \mathbf{C}^{-1},$$

which completes the proof of Theorem 2.

Apart from the many objections that could be lodged from a purely statistical standpoint, it would appear from Theorem 2 that it could be dangerous to over-relax excessively, that is, to employ values of ω close to 2. It is well-known, however Young [16], that, for certain problems, the optimal value of $\omega = \omega_b$, considered purely from the standpoint of the asymptotic rate of convergence, tends to 2 as $n \rightarrow \infty$. Furthermore, for ill-conditioned problems in general, where $\mu = \max_i |\mu_i|$ is very close to unity, then ω_b is very close to 2. The Young theory of successive over-relaxation also indicates that it is preferable to over-estimate ω_b rather than to underestimate it. It is in the context of these results that perhaps Theorem 2 may have some significance. Much depends upon the size of $\sigma^2 \mathbf{R}$, however, about which it is difficult to make a priori assumptions, and which will, in general, be a function of ω . The presence of \mathbf{C}^{-1} in (2.13) and (3.4) is merely confirmation of the well-known fact that the round-off error depends very definitely upon the conditioning of \mathbf{C} [12, 15, 2] even in an iterative method.

Furthermore, as was stated in the introduction, the central limit theorem is not applicable to the distribution of $\mathbf{r}^{(k)}$ as $k \rightarrow \infty$, since it may easily be shown that $\mathbf{r}^{(k)}$ remains uniformly bounded for all k . However, using Lemma 3 of Golub [5] in

precisely the same manner as is applied there to the Richardson second-order method, we may use Theorems 1 and 2 to obtain probabilistic bounds:

THEOREM 3. *Under the same hypotheses as in Theorem 2, then for all $0 \leq \beta \leq 1$,*

$$(3.36) \quad P\{\|\mathbf{r}^{(\infty)} - \mathbf{r}\|^2 \leq [\sigma^2 \operatorname{tr}(\mathbf{R}\mathbf{C}^{-1})]/[\omega(2 - \omega)\beta]\} \geq 1 - \beta,$$

where P denotes the probability function, $\operatorname{tr}(\dots)$ the trace of a matrix, $\|\dots\|$ the Euclidean norm, and

$$(3.37) \quad \mathbf{r}^{(\infty)} = \lim_{k \rightarrow \infty} \mathbf{r}^{(k)}, \quad \mathbf{r} = \lim_{k \rightarrow \infty} \mathbf{r}_k = \frac{1}{\omega} \mathbf{C}^{-1} \boldsymbol{\epsilon}.$$

We remark that, if $\mathbf{R} = \mathbf{I}$, then we may bound $\operatorname{tr}(\mathbf{C}^{-1})$ in (3.36) by $n/(1 - \mu)$, where $\mu = \max_i |\mu_i|$.

4. Numerical Experiment. As a simple numerical example, we consider the one-dimensional Dirichlet problem

$$(4.1) \quad \begin{aligned} y'' &= f(x) \\ y(0) &= \beta_0; \quad y(1) = \beta_1 \end{aligned}$$

The system of linear equations which arises when this problem is discretised has the form:

$$(4.2) \quad \mathbf{C}\mathbf{y} = -h^2\mathbf{f} + \mathbf{g},$$

where \mathbf{C} is the $n \times n$ matrix

$$(4.3) \quad \mathbf{C} = \boldsymbol{\pi} \begin{bmatrix} 2 & -1 & 0 & \cdot & \cdot & 0 \\ -1 & 2 & -1 & \cdot & \cdot & \cdot \\ 0 & -1 & 2 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & \cdot & \cdot & -1 & 2 \end{bmatrix} \boldsymbol{\pi}^*$$

for some permutation matrix, $\boldsymbol{\pi}$, and

$$(4.4) \quad \begin{aligned} \mathbf{f} &= \boldsymbol{\pi}[f_1, \dots, f_n]^*, \\ \mathbf{g} &= \boldsymbol{\pi}[\beta_0, 0, \dots, 0, \beta_1]^*, \end{aligned}$$

where $f_i = f(x_i)$ and $x_i = ih$, $(n + 1)h = 1$. In order to apply our theory, $\boldsymbol{\pi}$ is such that \mathbf{C} has the σ_1 -ordering. Now it can be shown (Marcus [8]) that if $\mathbf{C}^{-1} = (\gamma_{ij})$, then

$$(4.5) \quad \gamma_{ij} = \begin{cases} i(n + 1 - j)/(n + 1), & \text{if } j \geq i \\ j(n + 1 - i)/(n + 1), & \text{if } j \leq i \end{cases}$$

and hence, letting $\gamma_i = \sum_{j=1}^n \gamma_{ij}$,

$$(4.6) \quad \gamma_i = i(n + 1 - i)/2.$$

Suppose in Theorem 1 that

$$(4.7) \quad \boldsymbol{\epsilon} = \boldsymbol{\pi}[\epsilon_1, \epsilon_2, \dots, \epsilon_n]^*, \quad \mathbf{r} = \boldsymbol{\pi}[\rho_1, \rho_2, \dots, \rho_n]$$

Then, for some permutation $\pi = \pi(i)$ of $\{1, \dots, n\}$,

$$(4.8) \quad \rho_{\pi(i)} = \left\{ (n+1-i) \sum_{j < i} j \epsilon_j + i \sum_{j \geq i} (n+1-j) \epsilon_j \right\} / [\omega(n+1)]$$

If we assume that all the ϵ_j 's are constant, $\epsilon_j \equiv \epsilon$ (say), then

$$(4.9) \quad \rho_{\pi(i)} = \gamma_i \epsilon / \omega = i(n+1-i) \epsilon / 2\omega$$

Similarly for the variance, denoting, in Theorem 2, the diagonal elements of \mathbf{V} by v_i and assuming that $\mathbf{R} = \mathbf{I}$,

$$(4.10) \quad v_{\pi(i)} = \sigma^2 i(n+1-i) / [\omega(2-\omega)(n+1)]$$

In actual numerical computations, it is obviously difficult to make *a priori* assumptions about the size of ϵ and σ^2 , apart from justifying all the sundry assumptions that led to (5.9) and (5.10). In an attempt to facilitate these problems, we computed the solution to

$$(4.11) \quad \begin{aligned} y_m''(x) &= k_1(e^x + e^{1-x}) \\ y_m(0) &= y_m(1) = k_2 + m\Delta \end{aligned}$$

for certain constants k_1, k_2 and where Δ is a small increment. k_1 and k_2 were chosen ($k_1 = \frac{3}{16}, k_2 = \frac{5}{8}$) such that the analytical solution of (4.11) lay between $\frac{1}{2}$ and 1 for $m \leq 100$ and $\Delta = 10^{-4}$, and so that $y_m''(x) - y_m(x)$ was of small order of magnitude in comparison to $\frac{1}{2}$. The first of these requirements was to meet the demand that the ϵ_j 's be constant, and so that our floating-point arithmetic became effectively fixed-point, and the second of these requirements attempted to ensure that (2.5) was fulfilled, since then our initial guess, namely of setting y_i equal to f_i , was close to the exact solution and would not change by much during the computation.

The successive over-relaxation procedure was carried out in both single and double-precision arithmetic, the difference between the two being considered as the accumulated round-off error. The mean and variance of the latter were calculated, using $m = 1, \dots, 100$ as a sample (see Henrici [7] for a full discussion of this kind of experiment). One-hundred iterations were performed; it was surprisingly found that both the means and the variances converged, and that $k = 100$

TABLE 1
Dependence of Expected Values upon n

ω	1.2			1.5			1.8		
	(a)	(b)	(c)	(a)	(b)	(c)	(a)	(b)	(c)
10	61	71	67	64	71	56	65	71	57
20	247	257	212	253	257	183	251	257	246
30	564	560	517	576	560	432	572	560	478
40	1011	980	917	1064	980	820	1079	980	1053
50	1497	1517	1694	1504	1517	1596	1513	1517	1581

- (a) Experimental values of $10^{10} \rho_{\max}$, fixed-point.
 (b) Theoretical values of $10^{10} \rho_{\max}$, fixed-point.
 (c) Experimental values of $10^9 \rho_{\max}$, floating-point.

TABLE 2
Dependence of Variances upon n

ω	1.2			1.5			1.8		
n	(d)	(e)	(f)	(d)	(e)	(f)	(d)	(e)	(f)
10	131	76	117	137	121	144	462	320	381
20	247	144	206	342	230	315	1038	611	735
30	372	213	317	398	340	498	1316	903	1076
40	516	282	346	621	449	881	1521	1194	1298
50	724	350	590	740	559	1113	1976	1486	1727

(d) Experimental values of $10^{20} V_{\max}$, fixed-point.

(e) Theoretical values of $10^{20} V_{\max}$, fixed-point.

(f) Experimental values of $10^{18} V_{\max}$, floating-point.

TABLE 3
Dependence of Expected Values and Variances upon ω

ω	(a)	(b)	(c)	(d)	(e)	(f)
1.1	239	} 257	161	267	129	167
1.2	247		212	247	144	206
1.3	252		229	258	172	218
1.4	252		252	300	191	229
1.5	253		183	342	230	315
1.6	253		263	491	291	358
1.7	254		263	535	383	505
1.8	251		246	1038	611	735
1.82	274		217	1188	666	921
1.84	247		225	1201	752	1032
1.86	254		294	1441	844	1409
1.88	258		197	1539	988	1906
1.90	255		219	2005	1216	2124
1.92	254		240	2392	1525	2706
1.94	254		261	3335	2041	3120
1.96	274		270	4621	3073	4914
1.98	261		221	8742	6158	7091

was sufficient to ensure this convergence, which is required by the asymptotic nature of theorems 1 and 2. The procedure was carried out twice, once in fixed-point† and once in floating-point arithmetic, in both cases symmetric rounding being employed. An analysis of the arithmetic operations involved, using in particular Theorem 1.9 and equations (1-107) and (1-108) of Henrici [7], leads to the assumption that for the fixed-point computation

$$(4.12) \quad \epsilon \cong \frac{\omega u}{2}; \quad \sigma^2 \cong \frac{9\omega^2 + 16}{8} \frac{u^2}{12}$$

for $1 < \omega < 2$, where u is the basic machine unit. The fact that ϵ is not zero in

† The author would like to acknowledge the assistance of Mr. Peter Golitzen in this connection.

spite of symmetric rounding being employed stems from the particular equation solved, where we have to form products of the form $\frac{1}{2}y_i$ where $\frac{1}{2} < y_i \leq 1$, and the only possible values of the round-off error in forming this product is either 0 or $u/2$. This is also responsible for the peculiar expression for σ^2 .

Since the floating-point computations performed were in fact pseudo fixed-point, we also assume that (4.12) holds in this case, too. In the fixed-point calculation, $u = 2^{-30}$, whilst for the floating-point calculation, $u = 2^{-27}$. Thus from (4.9), (4.10) and (4.12), we expect that

$$(4.13) \quad \begin{aligned} \rho_{\pi(i)} &\cong \frac{1}{2}i(n+1-i)u/2 \\ v_{\pi(i)} &\cong [i(n+1-i)(9\omega^2+16)u^2]/[96(n+1)\omega(2-\omega)] \end{aligned}$$

so that if we let $\rho_{\max} = \max_i |\rho_i|$, $V_{\max} = \max_i |V_i|$,

$$(4.14) \quad \begin{aligned} \rho_{\max} &\cong (n+1)^2u/16 \\ V_{\max} &\cong [(n+1)(9\omega^2+16)u^2]/[384\omega(2-\omega)] \end{aligned}$$

Or computational results are divided into two parts, one to show dependence upon n , and the other to show dependence upon ω . In Tables 1 and 2 we show dependence upon n for $\omega = 1.2, 1.5; 1.8$. In Table 3 the dependence upon ω for $n = 20$ is shown; in this case ω_b , the optimal value of ω , is approximately 1.74. We give theoretical values only for the fixed-point computations. It will be seen that the experimentally calculated expected values correlate quite closely with the predicted theory, whereas the experimental values of the variances seem to be slightly higher than the theoretical values. This latter situation is most probably due to the non-independence of the local round-off errors after a large number of iterations. Similar experiences were obtained in the numerical experiments of Golub and Moore [6].

The author would like to extend his appreciation to Professor P. K. Henrici for his guidance in the preparation of the dissertation upon which this paper is principally based. He would also like to acknowledge the constructive criticisms of Dr. J. H. Wilkinson.

National Physical Laboratory
Teddington, Middlesex, England

1. A. A. ABRAMOV, "On the influence of round-off errors in the solution of Laplace's equation," *Vycis. Matem. i Vycis. Teh.*, v. 1, 1953, p. 37-41.
2. J. DESCLoux, "Note on the round-off error in iterative processes," *Math. Comp.*, v. 17, 1963, p. 18-27.
3. G. E. FORSYTHE, "Note on rounding-off errors," *SIAM Rev.*, v. 1, 1959, p. 66-67.
4. G. H. GOLUB, "The use of Chebyshev matrix polynomials in the iterative solution of linear equations compared to the method of successive over-relaxation," Doctoral Thesis, University of Illinois, 1959.
5. G. H. GOLUB, *Bounds for the Round-Off Errors in the Richardson Second-Order Method*, Nordisk Tidsskrift for Informations-Behandling, v. 2, 1962, p. 212-223.
6. G. H. GOLUB, & J. K. MOORE, *ibid* (appendix).
7. P. K. HENRICI, *Discrete-variable Methods in Ordinary Differential Equations*, John Wiley & Sons, Inc., New York, 1961.
8. M. MARCUS, "Basic theorems in matrix theory," *N.B.S. Appl. Math. Ser.* No. 57, 1960.
9. A. OSTROWSKI, "On the linear iteration procedures for symmetric matrices," *Rendic. di Mat. e.d.s. Applicaz.*, v. 13, 1954, p. 1-24.
10. W. SIBAGAKI, "On the idea of 'numerical convergence' and its applications," *Mem. Fac. Sci. Kyushu Univ. Ser. A*, v. 5, 1950, p. 89-97.

11. A. M. TURING, "Rounding errors in algebraic processes," *Quart. J. Appl. Math.*, v. 1, 1948, p. 287-307.
12. M. URABE, "Convergence of numerical iteration in solution of equations," *J. Sci. Hiroshima Univ. Ser. A*, v. 19, 1956, p. 479-489.
13. R. S. VARGA, *Matrix Iterative Analysis*, Prentice-Hall, New Jersey, 1962.
14. J. H. WILKINSON, "Rounding errors in algebraic processes," *Proc. Int. Conference on Information Processing*, UNESCO, 1959, p. 44.
15. J. H. WILKINSON, "Error analysis of direct methods of matrix inversion," *J. Assoc. Comp. Mach.*, v. 8, 1961, p. 281-330.
16. D. M. YOUNG, "Iterative methods for solving-partial differential equations of the elliptic type," *Trans. Amer. Math. Soc.*, v. 76, 1954, p. 92-111.