

Then equation (33), with $r \equiv 0$, may be shown to imply

$$(47) \quad (1 + q\Delta t/\Delta x)P_R^2 + (1 - q\Delta t/\Delta x)P_L^2 = (1 + q\Delta t/\Delta x)M_L^2 + (1 - q\Delta t/\Delta x)M_R^2 - 4p\sigma(1 + 2p\sigma)^{-2}\{M_L + M_R + (M_L - M_R)q\Delta t/\Delta x\}^2.$$

A front-function which displays a decreasing tendency may thus be written as

$$(48) \quad H = \sum_n (1 \pm q\Delta t/\Delta x)(\phi_{n,j(n)} - \phi_{n-1,j(n-1)})^2;$$

the \pm as $j(n-1) = j(n) \pm 1$.

In (47) q is properly written q_{nj} . In (48) it becomes ambiguous, being either $q_{n-1,j(n)}$ or $q_{n,j(n) \pm 1}$. Thus we can only assert that H is non-decreasing except by reason of variations in q with the advance of the front or disturbance introduced at the boundaries.

In equation (33) with nonvanishing r , the precise meaning of "stability" is not evident and no demonstration of properties approximating this concept are known to us. A preliminary qualitative examination failed to disclose any indication that variations in ϕ can grow to an alarming extent.

California Institute of Technology
and Continental Oil Company

E. C. DU FORT

California Institute of Technology

S. P. FRANKEL

We are grateful for the support and encouragement given this investigation by the Continental Oil Company. We have been aided by numerous discussions with L. BREIMAN.

¹ H. LEWY, "On the convergence of solutions of difference equations," *Studies and Essays Presented to R. Courant on his 60th Birthday*. New York, 1948, p. 211-214.

G. G. O'BRIEN, M. A. HYMAN, & S. A. KAPLAN, "A study of the numerical solution of partial differential equations," *Jn. Math. Phys.*, v. 29, 1951, p. 223-251.

F. JOHN, "On Integration of Parabolic Equations by Difference Methods," *Comm. Pure Appl. Math.*, v. 5, 1952, p. 155-211.

R. P. EDDY, *Stability in the Numerical Solution of Initial Value Problems in Partial differential Equations*. Naval Ordnance Lab., Memorandum 10232, 1949.

² S. P. FRANKEL, "Convergence rates of iterative treatments of partial differential equations," *MTAC*, v. 4, 1950, p. 65-75.

³ H. S. CARSLAW & J. C. JAEGER, *Conduction of Heat in Solids*. Oxford, 1947, p. 280-281.

On Over and Under Relaxation in the Theory of the Cyclic Single Step Iteration

In order to speed up the sometimes very tedious computations in solving equations by the single step method, the device of the so-called "incomplete relaxation" (under or over relaxation) has been often used, although apparently no systematic discussion of this device has been as yet tried.¹

It may seem, however, that in the case of the "relaxation procedure" the speeding up can be achieved in this way only in special cases, at least in the case of a symmetric positive definite matrix. Indeed, if the progress of the computation is measured by the decrease of a corresponding quadratic form $A(\zeta_k)$ depending on the k -th approximating vector ζ_k , we have the formula

$$(1) \quad A(\zeta_k) - A(\zeta_{k+1}) = q_k(2 - q_k) |r_{Nk}^{(k)}|^2 / a_{NkNk}$$

where N_k is the index of the variable the value of which is improved at the k -th step, $r_{N_k}^{(k)}$ the corresponding "residual", q_k a coefficient characterizing the degree of the incomplete relaxation at this step. (We have usually $0 < q_k \leq 2$; for $0 < q_k < 1$, we have *under relaxation* and for $1 < q_k \leq 2$ *over relaxation*, while for $q_k = 1$ we have the usual complete relaxation.)² However, the above formula shows that *the decrease of $A(\xi_k)$ is maximum* if q_k is taken as 1. Thus far the improvement in the case of $q_k \neq 1$ appears to be only possible if for such a value of q_k one of the residuals *at the next steps* will come out particularly large.

In what follows we discuss completely the case of a real 2×2 matrix, symmetric or not. It then turns out that, if the *cyclic* single step iteration converges, this convergence can indeed in all cases be essentially improved in using the incomplete relaxation with a convenient coefficient q , the same at each step. It turns out that even in many cases of unsymmetric matrices where the usual cyclic single step iteration diverges, it is made convergent in using a convenient incomplete relaxation.

The usual cyclic single step iteration is applied to a system of equations

$$(2) \quad \sum_{\nu=1}^n a_{\mu\nu} x_\nu = y_\mu \quad (\mu = 1, \dots, n)$$

with a non-singular matrix $A = (a_{\mu\nu})$. We can assume without loss of generality that $y_\mu = 0$ ($\mu = 1, \dots, n$), since we can always obtain this by changing the origin. Then we get from one approximation vector $\xi = (x_1, \dots, x_n)$ the next one $\xi' = (x_1', \dots, x_n')$ by the formula

$$(3) \quad a_{\mu\mu} x_\mu' = a_{\mu\mu} x_\mu - \sum_{\nu=\mu}^n a_{\mu\nu} x_\nu - \sum_{\nu=1}^{\mu-1} a_{\mu\nu} x_\nu' \quad (\mu = 1, \dots, n).$$

This is replaced for the incomplete relaxation with a coefficient q by

$$(4) \quad a_{\mu\mu} x_\mu' = a_{\mu\mu} x_\mu - q \sum_{\nu=\mu}^n a_{\mu\nu} x_\nu - q \sum_{\nu=1}^{\mu-1} a_{\mu\nu} x_\nu' \quad (\mu = 1, \dots, n)$$

and this can be written in the form

$$(5) \quad a_{\mu\mu} x_\mu' + q \sum_{\nu=1}^{\mu-1} a_{\mu\nu} x_\nu' = a_{\mu\mu} (1 - q)x_\mu - q \sum_{\nu=\mu+1}^n a_{\mu\nu} x_\nu.$$

We now decompose A into the sum

$$(6) \quad A = L + D + R$$

where D is the diagonal matrix containing the main diagonal of A , while in L all elements on the diagonal and to the right of it and in R all elements on the diagonal and to the left of it vanish. We make here the assumption, which is usually made in the theory of the single step iteration, that none of the diagonal elements of A vanishes; then we can write (5) in the form

$$(D + qL)\xi' = ((1 - q)D - qR)\xi$$

and in solving this, since the triangular matrix $D + qL$ is non-singular, we obtain

$$\xi' = (D + qL)^{-1} ((1 - q)D - qR)\xi.$$

If now we put

$$(7) \quad Q_q = (D + qL)^{-1} ((1 - q)D - qR)$$

we see that the k -th iterated vector ξ_k is obtained from the starting vector ξ_0 by the formula

$$\xi_k = Q_q^k \xi_0 \quad (k = 1, 2, \dots).$$

In order that our iteration be convergent for any starting vector ξ_0 , it is necessary and sufficient that the maximum modulus Λ_q of the *characteristics* roots of the matrix Q_q be less than 1. The speed of the convergence is then measured by Λ_q . The smaller Λ_q , the faster is the convergence.

The characteristic equation of (7),

$$|\lambda E - (D + qL)^{-1} ((1 - q)D - qR)| = 0,$$

becomes, if the matrix of the left hand expression is multiplied by $D + qL$,

$$(8) \quad |\lambda(D + qL) - ((1 - q)D - qR)| = 0.$$

We discuss this equation only in the case of a *real matrix* A with $n = 2$. Here the number

$$(9) \quad u = (a_{12}a_{21})/(a_{11}a_{22})$$

is the characteristic constant of the problem. Our results are contained in the following four theorems.

THEOREM 1. *If $u > 1$ then for all q from $(0, 2)$ we have $\Lambda_q > 1$ and the process is divergent.*

THEOREM 2. *Suppose that $u < 1$ and put*

$$(10) \quad q_0 = 2/(1 + (1 - u)^{\frac{1}{2}}).$$

Then with monotonically increasing q , Λ_q is monotonically decreasing for $q < q_0$ and monotonically increasing for $q > q_0$, so that the optimal value of Λ_q is

$$(11) \quad \Lambda_{opt} = \Lambda_{q_0} = |u|/(1 + (1 - u)^{\frac{1}{2}})^2 \quad (u < 1).$$

THEOREM 3. *Suppose that $u < 0$. Then a necessary and sufficient condition for the convergence of our procedure is that*

$$(12) \quad 0 < q < q_1 = 2(1 + |u|^{\frac{1}{2}})^{-1}.$$

If $0 < u < 1$, we have convergence for all q with $0 < q < 2$.

THEOREM 4. *Suppose that $|u| < 1$. Then a necessary and sufficient condition for $\Lambda_q < \Lambda_1$ is that q be contained in the interior of the interval between 1 and $1 + u$.*

To prove these theorems we start from the corresponding case of (8)

$$\begin{vmatrix} (\lambda + q - 1)a_{11} & qa_{12} \\ \lambda qa_{21} & (\lambda + q - 1)a_{22} \end{vmatrix} = 0.$$

Without loss of generality we can assume that $a_{11} > 0$, $a_{22} > 0$. In dividing the first row and the first column by $a_{11}^{\frac{1}{2}}$ and the second row and second column by $a_{22}^{\frac{1}{2}}$ we can reduce A to the form

$$A = \begin{pmatrix} 1 & \beta \\ \alpha & 1 \end{pmatrix}, \quad u = \alpha\beta.$$

Our equation for λ becomes now

$$(13) \quad N_q(\lambda) = \begin{vmatrix} \lambda + q - 1 & q\beta \\ \lambda q\alpha & \lambda + q - 1 \end{vmatrix} \\ = \lambda^2 + \lambda(2q - 2 - q^2u) + (q - 1)^2 = 0.$$

For $q = 1$ we obtain

$$(14) \quad \Lambda_1 = |u|$$

and we see that the usual cyclic single step iteration converges, for an arbitrary starting vector, only if $|u| < 1$.

In the following discussion we assume $u \neq 0$. The discriminant Δ of the polynomial $N_q(\lambda)$ is

$$\Delta = \frac{1}{4} q^2 u (uq^2 - 4q + 4).$$

This vanishes for $q = q_0 = 2/(1 + (1 - u)^{\frac{1}{2}})$, where obviously

$$(15) \quad \begin{cases} 1 < q_0 < 2 & (0 < u < 1) \\ 0 < q_0 < 1 & (u < 0), \end{cases}$$

while the other root $2/(1 - (1 - u)^{\frac{1}{2}})$ exceeds 2 if $0 < u < 1$ and is negative if $u < 0$. We see therefore that

$$(16) \quad \begin{cases} \Delta > 0 & (u > 1) \\ \text{Sgn } \Delta = \text{Sgn } u(q_0 - q) & (u < 1). \end{cases}$$

We consider now two cases according as $\Delta \leq 0$ or $\Delta > 0$. By (16), Δ is negative only if $u < 1$ and either $0 < u < 1$ and $q_0 < q \leq 2$ or $u < 0$ and $0 < q < q_0$. In both cases we have obviously from (13), $\Lambda_q = |q - 1|$ ($\Delta < 0$) and in particular

$$(17) \quad \begin{cases} \Lambda_q = q - 1 & (u > 0, \quad 2 \geq q \geq q_0) \\ \Lambda_q = 1 - q & (u < 0, \quad 0 < q \leq q_0). \end{cases}$$

In virtue of (16), the hypothesis $u > 1$ of theorem 1 is never realized for a negative Δ . Theorem 2 follows immediately from formulas (17). Theorem 3 follows from the fact that by (17) $\Lambda_q < 1$, while for $u < 0$, in any case $q_0 < q_1$ and q cannot exceed q_0 by (17). It follows finally from (17) that for $u > 0$ the condition

$$(18) \quad \Lambda_q < \Lambda_1 = |u|$$

is equivalent to $q < 1 + u$ while q remains $\geq q_0 > 1$. On the other hand, if we have $u < 0$, the condition (18) is equivalent to $q > 1 + u$, while q remains $\leq q_0 < 1$. Therefore theorem 4 and all our assertions are true in case $\Delta \leq 0$.

Now let $\Delta > 0$. If λ is the root of $N_q(\lambda)$ with $|\lambda| = \Lambda_q$, we have

$$(19) \quad \lambda = R + \epsilon\Delta^{\frac{1}{2}}, \quad R = \frac{u}{2} \left(q^2 - 2 \frac{q-1}{u} \right),$$

where $\epsilon = \text{Sgn } R$. We have obviously

$$(20) \quad \Lambda_q = \epsilon\lambda.$$

The monotonicity of λ with respect to q depends on the sign of

$$\frac{d\lambda}{dq} = 2 \frac{qu\lambda - \lambda - (q-1)}{2\lambda - q^2u - 2 + 2q}.$$

By (19), the denominator is equal to

$$2(\lambda - R) = 2\epsilon\Delta^\dagger$$

and has the sign of ϵ . If we denote the numerator by δ , we have

$$(21) \quad \delta = \lambda(qu - 1) + 1 - q$$

and therefore, by (20),

$$(22) \quad \text{Sgn} \frac{d\Lambda_q}{dq} = \text{Sgn} \delta.$$

We deal first with the case $u > 1$. Here we have $\Delta > 0$, the roots of (13) are real and, since

$$N_q(1) = 1 - q^2u - 2 + 2q + q^2 - 2q + 1 = q^2(1 - u) < 0,$$

one root λ exceeds 1, and so $\Lambda_q > 1$. We see that in this case we have divergence for any value of $q > 0$ and theorem 1 is proved.

Since for $u = 1$, A becomes singular, from now on we can make the assumptions

$$(23) \quad u < 1, \quad q_0 \text{ is real, } \Delta > 0.$$

We have then in particular from (16)

$$(24) \quad \begin{cases} \text{either } u > 0, & q < q_0, & q_0 > 1 \\ \text{or } u < 0, & q > q_0, & q_0 < 1. \end{cases}$$

We prove now the following lemma, the proof of which is the main difficulty of the paper:

LEMMA. *Under the assumptions (23) $u\delta$ is always negative.*

Denote by δ_1 and δ_2 the two values of δ corresponding to the two roots of $N_q(\lambda)$. We have from (21) and (13) after some simplifications

$$(25) \quad \delta_1 + \delta_2 = qu(q^2u - 3q + 2),$$

$$(26) \quad \delta_1\delta_2 = (q-1)uq^2(1-u).$$

The expression on the right in (25) vanishes for $q = q_2$, where

$$(27) \quad q_2 = 4/(3 + (9 - 8u)^\dagger) < 1 \quad (u < 1),$$

while the other root exceeds 2 for $u > 0$ and is negative for $u < 0$. Therefore we have

$$(28) \quad \text{Sgn}(\delta_1 + \delta_2) = \text{Sgn} u(q_2 - q).$$

We will now consider separately the cases $u > 0$ and $u < 0$, and prove in the first case $\delta < 0$ and in the second $\delta > 0$.

In the case $u > 0$ we have $1 > u > 0$, and, since $\Delta > 0$, $q < q_0$. From (24) and (27) it follows that

$$(29) \quad 2 > q_0 > 1 > q_2.$$

If q exceeds 1, it follows from (29), (28), and (26) that

$$\delta_1 + \delta_2 < 0, \quad \delta_1\delta_2 > 0, \quad \delta_1 < 0, \quad \delta_2 < 0$$

and therefore $\delta < 0$.

If $q < 1$, it follows from (26) that $\delta_1\delta_2 < 0$. On the other hand the expression of R in (19), for $u > 0$, $q < 1$, becomes positive. We have therefore in this case $\epsilon = 1$, λ is the greater root of $N_q(\lambda)$ and since $qu - 1$ in (21) becomes negative, we have

$$\delta = \text{Min}(\delta_1, \delta_2).$$

But then it follows from $\delta_1\delta_2 < 0$ that δ is negative in this case also.

We consider now the case $u < 0$. Since Δ is assumed positive, we have here by (16), $q > q_0$, while on the other hand from (24) and (27) it follows that

$$q_2 < q_0 < 1.$$

Since therefore in any case $q > q_2$, we have from (28)

$$(30) \quad \delta_1 + \delta_2 > 0.$$

On the other hand we conclude from (26)

$$(31) \quad \text{Sgn } \delta_1\delta_2 = \text{Sgn}(1 - q).$$

If therefore $q < 1$, we have $\delta_1\delta_2 > 0$ and from (30)

$$\delta_1 > 0, \quad \delta_2 > 0, \quad \delta > 0.$$

Suppose now $q > 1$; since in this case, by (19), R is negative, we have $\epsilon = -1$, $\Lambda_q = -\lambda$ and from (21) it follows now, since the coefficient of λ in (21) is negative, that

$$\delta = \text{Max}(\delta_1, \delta_2).$$

But now we see from (30) that δ is again positive. Hence our lemma is proved.

We see now, from (22), that in the case of positive Δ , Λ_q monotonically decreases for $q < q_0$ and monotonically increases for $q > q_0$, as we already deduced from (17) in the case of negative Δ .

The minimum of Λ_q is obtained for $q = q_0$. We have, in applying for instance (17) and in using (10),

$$(32) \quad \Lambda_{q_0} = |(1 - (1 - u)^{\frac{1}{2}})/(1 + (1 - u)^{\frac{1}{2}})| \\ = |u|(1 + (1 - u)^{\frac{1}{2}})^{-2} \quad (u < 1),$$

and theorem 2 is completely proved.

The expression (32) is less than $|u|$ unless $u = 1$. For $|u| < 1$, we therefore always have an improvement in using incomplete relaxation with $q = q_0$, which is particularly pronounced for small $|u|$, since we have

$$(33) \quad \Lambda_{q_0} \sim |u|/4 \quad (u \rightarrow 0).$$

As to the values of q for which we have convergence at all, we have, since $\Lambda_0 = 1$, convergence when q is in the interval $(0, q_0)$. On the other hand, the roots of (13) become, for $q = 2$,

$$(34) \quad 2u - 1 \pm 2(u^2 - u)^{\frac{1}{2}};$$

they are complex for $1 > u > 0$ and it follows from (17) that

$$(35) \quad \Lambda_2 = 1 \quad (u > 0).$$

For $1 > u > 0$ we therefore have convergence for all q from $(0, 2)$.

If, on the contrary, $u < 0$, we obtain

$$\Lambda_2 = 1 - 2u + 2(u^2 - u)^{\frac{1}{2}} \quad (u < 0)$$

and this exceeds 1. To obtain the value q_1 of q between q_0 and 2 for which Λ_q becomes 1, observe that

$$N_q(-1) = (2 - q)^2 + uq^2, \quad N_q(1) = q^2(1 - u).$$

For $u < 0$, $N_q(1)$ does not vanish while $N_q(-1)$ vanishes for

$$(36) \quad q_1 = 2/(1 + |u|^{\frac{1}{2}}) \quad (u < 0).$$

Thus q_1 lies always between q_0 and 2, and since for $q = q_1$ the product of two roots of (13) is $(q_1 - 1)^2 < 1$, we have indeed

$$\Lambda_{q_1} = 1 \quad (u < 0).$$

We have therefore for $u < 0$ convergence if and only if q lies in the interval $(0, q_1)$. In particular we have here for suitable values of q convergence for all negative u , but these values become small for large value of $-u$. Theorem 3 is thus completely proved.

It remains finally to answer the question for what values of q does the incomplete relaxation give any improvement at all, that is to say, $\Lambda_q < |u|$ (of course we assume here $|u| < 1$).

Now we have, by (14), $\Lambda_1 = |u|$ and the same is true for $q = 1 + u$:

$$(37) \quad \Lambda_1 = |u| = \Lambda_{1+u} \quad (|u| < 1).$$

Indeed, we verify immediately that $1 + u \geq q_0$ according as u is positive or negative; therefore the roots of (13) are complex for $q = 1 + u$, and (37) follows from (17).

But now it follows from theorem 2 that we have improvement in the case of the incomplete relaxation (in contrast to the case $q = 1$) if and only if q lies in the interior of the interval between 1 and $1 + u$. Theorem 4 is thus proved.

We discuss finally a special example of a symmetric 2×2 matrix in which the improvement of the convergence due to the introduction of the incomplete relaxation is easily demonstrated explicitly.

$$\text{Let} \quad A = \begin{pmatrix} 1 & 3/5 \\ 3/5 & 1 \end{pmatrix}.$$

We have in this case by (10), (11), and (14)

$$(38) \quad \Lambda_1 = u = 9/25 > 1/3, \quad q_0 = 10/9, \quad \Lambda_{q_0} = 1/9.$$

The formulas for the cyclic single step iteration with $q = 1$ are in this case

$$(39) \quad x_1^{(r+1)} = -3x_2^{(r)}/5, \quad x_2^{(r+1)} = -3x_1^{(r)}/5$$

and it is readily verified that, in putting $\bar{\xi}_r = (x_1^{(r)}, x_2^{(r)})$, we have

$$\bar{\xi}_r = (9/25)^r \bar{\xi}_0.$$

In taking $\bar{\xi}_0 = (1, 1)$ we have then

$$(40) \quad |\bar{\xi}_\nu| = \sqrt{2}(9/25)^\nu.$$

Consider on the other hand the over relaxation with the value of $q = 10/9$. Here the components of the approximating vectors are to be computed from the equations

$$x_1^{(\nu+1)} = -x_1^{(\nu)}/9 - 2x_2^{(\nu)}/3, \quad x_2^{(\nu+1)} = -x_2^{(\nu)}/9 - 2x_1^{(\nu+1)}/3.$$

If we put $\xi_\nu = (x_1^{(\nu)}, x_2^{(\nu)})$ and assume again $\xi_0 = (1, 1) = \bar{\xi}_0$, we obtain as is readily verified

$$\xi_\nu = 3^{-2\nu-1}(3 - 24\nu, 3 + 8\nu) \quad (\nu = 0, 1, \dots).$$

Here we have

$$|\xi_\nu| \sim 8(10)^{\frac{1}{2}\nu}9^{-\nu}/3 \quad (\nu \rightarrow \infty).$$

We give in what follows a table of the initial values of $|\bar{\xi}_\nu|$ and $|\xi_\nu|$.

ν	$ \bar{\xi}_\nu $	$ \xi_\nu $
1	0.5091	0.8780
2	0.1833	0.2010
3	0.06598	0.03388
4	0.02375	0.005048
5	0.008551	0.0007037

Although the difference between q_0 and 1 is very small, in fact $1/9$, the improvement is already observed at ξ_3 and becomes more and more pronounced from there on.

A. OSTROWSKI

American University, Washington, D. C.
University of Basle, Switzerland

This paper was prepared under a contract of the National Bureau of Standards with the American University, Washington, D. C.

¹S. P. FRANKEL, "Convergence rates of iterative treatments of partial differential equations," *MTAC*, v. 4, 1950, p. 65-75.

²A. M. OSTROWSKI, *On the Linear Iteration Procedures for Symmetric Matrices*. NBS Report no. 1844, August 1952, p. 23, (68).

The Accuracy of Numerical Solutions of Ordinary Differential Equations

1. Introduction. The present paper describes a general method by which the random and systematic errors may be estimated of numerical solutions of any systems of ordinary differential equations. The errors arise from the accumulation of rounding-off errors, and from the use of erroneous formulas for performing the numerical integrations. The estimation is based on the properties of the solutions of the system of equations adjoint to the variational equations of the problem, and is applicable to any method of integration.