

A new method for the solution of eigenvalue problems

By M. R. Osborne*

One among the class of iterations proposed by Osborne and Michaelson (1964) is discussed. It is shown that for the ordinary eigenvalue problem this iteration is of third order, and that it may readily be modified to be of third order for the general eigenvalue problem provided the eigenvalue parameter occurs linearly. The analysis also suggests a further modification which gives a third-order iteration for the non-linear problem.

1. Introduction

In a recent paper (Osborne and Michaelson (1964)) the author has discussed the solution of eigenvalue problems having the form

$$M(\lambda)v = 0. \quad (1.1)$$

Here the eigenvalues are found as the zeros of a function $\beta(\lambda)$ defined by imposing a scaling condition on the solution vector v in the equation

$$M(\lambda)v = \beta(\lambda)x. \quad (1.2)$$

By scaling v so that one component is fixed in value independent of λ , and by applying Newton's method to the resulting function $\beta(\lambda)$, Osborne and Michaelson derive the iteration

$$\left. \begin{aligned} M(\lambda_i)v_{i+1} &= x_i/(x_i)_{p_i}, \\ M(\lambda_i)x_{i+1} &= \frac{dM}{d\lambda}(\lambda_i)v_{i+1} \\ \lambda_{i+1} &= \lambda_i - \frac{(v_{i+1})_{p_{i+1}}}{(x_{i+1})_{p_{i+1}}} \end{aligned} \right\} \quad (1.3)$$

where p_i is the index of the component of maximum modulus in x_i . In this paper iterations similar to (1.3) are discussed with a view to determining their order. (We say an iteration is of order p if δ_{j+1} is proportional to δ_j^p where δ_j is a parameter giving a measure of the error at the j th stage of the iteration.)

In the next section the ordinary eigenvalue problem is discussed ($\frac{dM}{d\lambda} = -I$). In this case the iteration of Osborne and Michaelson is of third order. Conditions sufficient to ensure the convergence of the iteration are also given. In Section 3 the case $\frac{dM}{d\lambda} = -B$, where B is a constant matrix is discussed. A modification of the basic iteration (1.3) is introduced and shown to be of third order. The iteration (1.3) is second order in this case. In Section 4, the case where M depends nonlinearly on λ is considered. Here both the iteration (1.3) and the modified iteration are of second order in general.

However, the analysis suggests how the iteration might be modified to give third-order convergence, and it also suggests a possible second-order iteration which would require the solution of only one set of linear equations per step.

In the case in which M depends linearly on λ the iteration (1.3) amounts to two stages of inverse iteration followed by a shift in the origin of λ . It is of considerable interest that a third-order process can be obtained merely by carrying out two inverse iterations, keeping the same λ (requiring one triangular factorization and two forward and back substitutions), as inverse iteration is anyway recommended for the calculation of eigenvectors (Wilkinson (1958)).

It must be stressed that the iterations discussed in this paper apply only for simple eigenvalues, and must be modified to deal with repeated eigenvalues.

2. The ordinary eigenvalue problem

In this section it is shown that the iteration

$$\left. \begin{aligned} [A - \lambda_i I]v_{i+1} &= x_i/(x_i)_{p_i}, \\ [A - \lambda_i I]x_{i+1} &= -v_{i+1}, \\ \lambda_{i+1} &= \lambda_i - \frac{(v_{i+1})_{p_{i+1}}}{(x_{i+1})_{p_{i+1}}} \end{aligned} \right\} \quad (2.1)$$

is of third order. It is assumed that the eigenvalue being determined is distinct. Let this eigenvalue be μ_1 , then the precise assumption made is that

$$\min_j (|\mu_j - \mu_1|, |\mu_j - \lambda_i|) \geq a \quad (2.2)$$

where the index j runs over the other eigenvalues of A . The need for the introduction of the parameter a , and its numerical significance, are discussed in Wilkinson (1961). Initially it is also assumed that the matrix A has a complete set of eigenvectors u_i scaled so that their maximum norm is 1. Whenever norms are used in this paper the maximum norm is assumed.

Further, it is assumed that $x_i/(x_i)_{p_i}$ is an approximation to u_1 in the sense that

$$x_i/(x_i)_{p_i} = \alpha_1 u_1 + q_i, \quad (2.3)$$

* Edinburgh University Computer Unit, 7, Buccleuch Place, Edinburgh 8.

where
$$\mathbf{q}_i = \sum_{j=2}^n \alpha_j^i \mathbf{u}_j, \tag{2.4}$$

and
$$\|\mathbf{q}_i\| \leq \sum_{j=2}^n |\alpha_j^i| \leq \delta_i, \tag{2.5}$$

and that λ_i is an approximation to μ_1 in the sense that
$$\lambda_i = \mu_1 + \delta_i \eta_i. \tag{2.6}$$

It is assumed that the parameter δ is small, and that the parameter η is of the same order of magnitude as a .

Theorem. If $\delta_i < \frac{1}{3}$, and $|\eta_i| < 1 \cdot 5a$, then for any $j \geq i$

$$\delta_{j+1} < 5 \cdot 4 \delta_j^3,$$

$$|\eta_{j+1}| < 1 \cdot 5a,$$

and the iteration is convergent.

Proof. From equation (2.1) it follows that

$$\mathbf{v}_{i+1} = -\frac{\alpha_1^i}{\delta_i \eta_i} \mathbf{u}_1 + \mathbf{r}_{i+1} \tag{2.7}$$

where
$$\mathbf{r}_{i+1} = \sum_{j=2}^n \frac{\alpha_j^i}{\mu_j - \lambda_i} \mathbf{u}_j, \tag{2.8}$$

and
$$\|\mathbf{r}_{i+1}\| \leq \delta_i/a. \tag{2.9}$$

Also
$$\mathbf{x}_{i+1} = -\frac{\alpha_1^i}{(\delta_i \eta_i)^2} \mathbf{u}_1 + \mathbf{s}_{i+1}, \tag{2.10}$$

where
$$\mathbf{s}_{i+1} = -\sum_{j=2}^n \frac{\alpha_j^i}{(\mu_j - \lambda_i)^2} \mathbf{u}_j, \tag{2.11}$$

and
$$\|\mathbf{s}_{i+1}\| \leq \delta_i/a^2. \tag{2.12}$$

Now, by equations (2.1) and (2.6),

$$\begin{aligned} \mu_1 - \lambda_{i+1} &= \mu_1 - \lambda_i + \frac{\frac{\alpha_1^i}{\delta_i \eta_i} (\mathbf{u}_1)_{p_{i+1}} - (\mathbf{r}_{i+1})_{p_{i+1}}}{\frac{\alpha_1^i}{(\delta_i \eta_i)^2} (\mathbf{u}_1)_{p_{i+1}} - (\mathbf{s}_{i+1})_{p_{i+1}}} \\ &= \delta_i \eta_i \frac{(\delta_i \eta_i)^2 (\mathbf{s}_{i+1})_{p_{i+1}} - (\delta_i \eta_i) (\mathbf{r}_{i+1})_{p_{i+1}}}{\alpha_1^i (\mathbf{u}_1)_{p_{i+1}} - (\delta_i \eta_i)^2 (\mathbf{s}_{i+1})_{p_{i+1}}}. \end{aligned} \tag{2.13}$$

Also, by equations (2.3) and (2.5)

$$1 - \delta_i \leq |\alpha_1^i| \leq 1 + \delta_i,$$

and by equations (2.10) and (2.12),

$$\left. |\alpha_1^i| - \frac{\delta_i^3 \eta_i^2}{a^2} \leq |\alpha_1^i (\mathbf{u}_1)_{p_{i+1}}| + \frac{\delta_i^3 \eta_i^2}{a^2} \right\} \tag{2.14}$$

so that $|\alpha_1^i (\mathbf{u}_1)_{p_{i+1}}| \geq 1 - \delta_i - 2\delta_i^3 \eta_i^2/a^2$

(It is readily verified that the numerical values of δ_i and $|\eta_i/a|$ stated in the theorem permit these inequalities to be satisfied with positive values on both sides of the inequality).

By taking moduli in equation (2.13), and using equations (2.9), (2.12), and (2.14) there is obtained

$$\delta_{i+1} |\eta_{i+1}| \leq \frac{\delta_i^3 \eta_i^2}{a} \frac{1 + \delta_i |\eta_i/a|}{1 - \delta_i - 3\delta_i^3 \eta_i^2/a^2}. \tag{2.15}$$

Writing γ_i for $|\eta_i/a|$, and dividing both sides of equation (2.16) by a gives

$$\delta_{i+1} \gamma_{i+1} \leq \delta_i^3 \gamma_i^2 \frac{1 + \delta_i \gamma_i}{1 - \delta_i (1 + 3\gamma_i^2 \delta_i^2)}. \tag{2.16}$$

Further $\mathbf{x}_{i+1}/(\mathbf{x}_{i+1})_{p_{i+1}} = \alpha_1^{i+1} \mathbf{u}_1 + \mathbf{q}_{i+1}$

where, by equations (2.10), (2.12), and (2.14)

$$\|\mathbf{q}_{i+1}\| \leq \delta_{i+1} = \frac{\delta_i^3 \gamma_i^2}{1 - \delta_i (1 + 3\gamma_i^2 \delta_i^2)}. \tag{2.17}$$

(Clearly, by the manner of its construction, the right-hand side of equation (2.17) provides a bound for $\sum_{j=2}^n |\alpha_j^{i+1}|$ so that it can be used to define δ_{i+1}). From equations (2.16) and (2.17) it is seen that

$$\delta_{i+1} \gamma_{i+1} \leq \delta_{i+1} (1 + \delta_i \gamma_i) \tag{2.18}$$

whence

$$\gamma_{i+1} \leq 1 + \delta_i \gamma_i. \tag{2.18}$$

Substituting the numerical values of bounds for δ_i and γ_i (noting that $\gamma_i \leq 1 \cdot 5$ for all values of a) shows that

- (a) from equation (2.18), $\gamma_{i+1} \leq 1 \cdot 5$
- (b) from equation (2.17), $\delta_{i+1} \leq 1/5 < 1/3$.

Therefore the conditions of the theorem are satisfied for $j = i + 1$ (and hence for all $j > i$) if they hold for $j = i$. The first part of the theorem is now obtained by inserting bounds into the coefficient of δ_i^3 in equation (2.17), while the second part follows immediately from (a) above. A direct computation using equation (2.15) gives

$$\delta_{i+1} |\eta_{i+1}| \leq 0 \cdot 6 \delta_i |\eta_i| \tag{2.19}$$

and as the bounds used in obtaining this inequality are valid for all $j \geq i$, the inequality is also valid for all $j \geq i$. This, plus the result that $\delta_{j+1} < 5 \cdot 4 \delta_j^3$, shows that the iteration converges.

In the case where the matrix A has principal vectors of grades > 1 the argument given above requires some modification. The crucial steps in the argument involve the bounding of $\|\mathbf{r}_{i+1}\|$ and $\|\mathbf{s}_{i+1}\|$ (equations (2.9) and (2.12)). Assume that \mathbf{u}_p is a principal vector of grade two associated with the eigenvalue μ_r . Then

$$(A - \mu_r I)^2 \mathbf{u}_p = 0,$$

whence

$$(A - \lambda_i I) \mathbf{u}_p + 2(\mu_r - \lambda_i) \mathbf{u}_p + (\mu_r - \lambda_i)^2 (A - \lambda_i I)^{-1} \mathbf{u}_p = 0. \tag{2.20}$$

Taking norms in this equation gives

$$\|(A - \lambda_i I)^{-1} \mathbf{u}_p\| \leq \frac{2\|\mathbf{u}_p\|}{a} + \frac{\|A - \lambda_i I\| \|\mathbf{u}_p\|}{a^2} \tag{2.21}$$

and

$$\|(A - \lambda_i I)^{-2} \mathbf{u}_p\| \leq \frac{2}{a} \|(A - \lambda_i I)^{-1} \mathbf{u}_p\| + \frac{\|\mathbf{u}_p\|}{a^2}. \quad (2.22)$$

The most significant feature of equations (2.21) and (2.22) is the occurrence of terms in $1/a^2$ and $1/a^3$ respectively. It is readily verified that the corresponding inequalities for a principal vector of grade m involve inverse powers of a of orders up to the m th and $(m + 1)$ th respectively. A possible corollary of this is that the numerical difficulties usually experienced when a is necessarily small are intensified if the close eigenvalue has associated with it principal vectors of grades higher than one.

3. A more general eigenvalue problem

In this section is considered the eigenvalue problem

$$[A - \lambda B] \mathbf{v} = 0. \quad (3.1)$$

If the matrix B has an inverse then the problem stated in equation (3.1) can be reduced to that considered in Section 2 by multiplying by B^{-1} . The iteration (2.1) applied to the matrix $B^{-1}A$ is equivalent to the iteration

$$\left. \begin{aligned} [A - \lambda_i B] \mathbf{v}_{i+1} &= B \mathbf{x}_i / (\mathbf{x}_i)_{p_i} \\ [A - \lambda_i B] \mathbf{x}_{i+1} &= -B \mathbf{v}_{i+1} \\ \lambda_{i+1} &= \lambda_i - \frac{(\mathbf{v}_{i+1})_{p_{i+1}}}{(\mathbf{x}_{i+1})_{p_{i+1}}} \end{aligned} \right\}. \quad (3.2)$$

This iteration differs significantly from that proposed by Osborne and Michaelson in the definition of \mathbf{v}_{i+1} . It will be shown to be of third order even when the matrix B does not have an inverse. The iteration proposed by Osborne and Michaelson gives only second order convergence for the problem (3.1).

First the solution of equation (3.3) is considered

$$[A - \lambda_i B] \mathbf{y} = \mathbf{x} \quad (3.3)$$

where $\lambda_i = \lambda + \delta_i$, and λ is an eigenvalue of equation (3.1), so that 0 is an eigenvalue of the matrix $[A - \lambda B]$. It is assumed that the other eigenvalues of the matrix $[A - \lambda B]$ are μ_2, \dots, μ_n , that \mathbf{u}_1 is the eigenvector corresponding to the zero eigenvalue, and that the μ_i , $i = 2, \dots, n$, satisfy the condition (2.2). The similarity normal form of $[A - \lambda B]$ is written TMT^{-1} , where $M_{11} = 0$. The vectors \mathbf{y} and \mathbf{x} have representations which are written

$$\mathbf{y} = T\boldsymbol{\alpha}, \quad \mathbf{x} = T\boldsymbol{\beta}. \quad (3.4)$$

In this notation equation (3.3) becomes

$$\begin{aligned} [A - \lambda_i B] \mathbf{y} &= [A - \lambda B] \mathbf{y} - \delta_i B \mathbf{y} \\ &= T M \boldsymbol{\alpha} - \delta_i B T \boldsymbol{\alpha} \\ &= T \boldsymbol{\beta} \end{aligned}$$

so that, defining G to be the matrix $-T^{-1}BT$,

$$[M + \delta_i G] \boldsymbol{\alpha} = \boldsymbol{\beta}. \quad (3.5)$$

Equation (3.5) is represented in the partitioned form

$$\begin{bmatrix} \delta_i G_{11} & \delta_i \mathbf{g}^T \\ \delta_i \mathbf{f} & \mathbf{K} \end{bmatrix} \begin{bmatrix} \boldsymbol{\alpha}_1 \\ \boldsymbol{\xi} \end{bmatrix} = \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\eta} \end{bmatrix}. \quad (3.6)$$

If δ_i is small enough then (by the condition (2.2)) \mathbf{K} possesses an inverse, and equation (3.6) gives

$$\left. \begin{aligned} \boldsymbol{\xi} &= \mathbf{K}^{-1}(\boldsymbol{\eta} - \boldsymbol{\alpha}_1 \delta_i \mathbf{f}) \\ \delta_i G_{11} \boldsymbol{\alpha}_1 + \delta_i \mathbf{g}^T \mathbf{K}^{-1}(\boldsymbol{\eta} - \boldsymbol{\alpha}_1 \delta_i \mathbf{f}) &= \boldsymbol{\beta}_1 \end{aligned} \right\} \quad (3.7)$$

whence

$$\left. \begin{aligned} \boldsymbol{\alpha}_1 &= \frac{\boldsymbol{\beta}_1}{\delta_i G_{11}} + \left(\frac{\boldsymbol{\beta}_1}{G_{11}^2} \mathbf{g}^T \mathbf{K}^{-1} \mathbf{f} - \frac{\mathbf{g}^T \mathbf{K}^{-1} \boldsymbol{\eta}}{G_{11}} \right) + O(\delta_i), \\ \boldsymbol{\xi} &= \mathbf{K}^{-1} \boldsymbol{\eta} - \frac{\boldsymbol{\beta}_1}{G_{11}} \mathbf{f} + O(\delta_i). \end{aligned} \right\} \quad (3.8)$$

Notes.—(1) Because of the term $\frac{\boldsymbol{\beta}_1}{G_{11}} \mathbf{f}$, $\boldsymbol{\xi}$ can be $O(1)$ even when $\boldsymbol{\eta}$ is $O(\delta_i)$ so that the projection of the vector \mathbf{x} on vectors other than the appropriate eigenvector is small. It is this that causes the iteration of Osborne and Michaelson to be of second order only. In particular, for this iteration, equation (3.12) does not hold.

(2) It is the matrix \mathbf{K}^{-1} that contains the terms $O(1/a)$ when the spacing between the eigenvalues is small.

(3) If $G_{11} = 0$ then, by equation (3.7), $\boldsymbol{\alpha}_1 = O(1/\delta_i^2)$. In this case it is readily seen that the iteration (3.2) cannot be satisfactory if indeed it converges at all.

Now let $\mathbf{x}_i / (\mathbf{x}_i)_{p_i}$ be an approximation to \mathbf{u}_1 in the sense that

$$\mathbf{x}_i / (\mathbf{x}_i)_{p_i} = a_i \mathbf{u}_1 + \delta_i \mathbf{q}_i \quad (3.9)$$

where a_i and $\|\mathbf{q}_i\|$ are $O(1)$. From equation (3.1)

$$[A - (\lambda_i - \delta_i)B] (\mathbf{x}_i / (\mathbf{x}_i)_{p_i} - \delta_i \mathbf{q}_i) = 0$$

whence

$$[A - \lambda_i B] a_i \mathbf{u}_1 = -\delta_i B \mathbf{x}_i / (\mathbf{x}_i)_{p_i} + \delta_i^2 B \mathbf{q}_i. \quad (3.10)$$

Also, from equation (3.2),

$$[A - \lambda_i B] \delta_i \mathbf{v}_{i+1} = \delta_i B \mathbf{x}_i / (\mathbf{x}_i)_{p_i}$$

so that

$$[A - \lambda_i B] (\delta_i \mathbf{v}_{i+1} + a_i \mathbf{u}_1) = \delta_i^2 B \mathbf{q}_i. \quad (3.11)$$

Using equation (3.8) with $\boldsymbol{\beta} = T^{-1}B \mathbf{q}_i$ gives

$$\delta_i \mathbf{v}_{i+1} + a_i \mathbf{u}_1 = \left[\delta_i \frac{\boldsymbol{\beta}_1}{G_{11}} + O(\delta_i^2) \right] \mathbf{u}_1 + \mathbf{r}_{i+1} \quad (3.12)$$

where $\|\mathbf{r}_{i+1}\| = O(\delta_i^2)$.

$$\text{Thus } \mathbf{v}_{i+1} = a_i^* \mathbf{u}_1 + \delta_i \mathbf{r}_{i+1}^* \quad (3.13)$$

where $a_i^* = O(1/\delta_i)$, and $\|\mathbf{r}_{i+1}^*\| = O(1)$.

The argument used for deriving \mathbf{v}_{i+1} can be used again to derive \mathbf{x}_{i+1} . Defining $\boldsymbol{\gamma}$ to be $T^{-1}B \mathbf{r}_{i+1}^*$, then

$$-\delta_i \mathbf{x}_{i+1} + a_i^* \mathbf{u}_1 = \left[\delta_i \frac{\boldsymbol{\gamma}_1}{G_{11}} + O(\delta_i^2) \right] \mathbf{u}_1 + \mathbf{s}_{i+1} \quad (3.14)$$

where $\|s_{i+1}\| = O(\delta_i^2)$. (3.15)

Define s_{i+1}^* to be the vector $\delta_i^{-2}s_{i+1}$.

The error in the next approximation to the eigenvalue can now be estimated. From equation (3.2)

$$\begin{aligned} \lambda_{i+1} &= \lambda + \delta_{i+1} \\ &= \lambda + \delta_i - \frac{a_i^*(u_1)_{p_{i+1}} + \delta_i(v_{i+1}^*)_{p_{i+1}}}{\delta_i \left[1 - \frac{\delta_i}{a_i^*} \frac{\gamma_1}{G_{11}} + O(\delta_i^2) \right]} (u_1)_{p_{i+1}} - \delta_i(s_{i+1})_{p_{i+1}} \end{aligned}$$

where the fact that $a_i^* = O(1/\delta_i)$ has been used, $= \lambda + O(\delta_i^3)$. (3.16)

It follows from equations (3.14), (3.15) and (3.16) that the iteration (3.2) is of third order.

4. The non-linear problem

The general form of the eigenvalue problem is

$$M(\lambda)u_1 = 0. \quad (4.1)$$

Iterative methods for the solution of this problem are considered in this section. The notation of Section 3 is followed here.

Expanding equation (4.1) by Taylor's theorem gives

$$\left\{ M(\lambda_i) - \delta_i \frac{dM}{d\lambda}(\lambda_i) + \frac{\delta_i^2}{2} \frac{d^2M}{d\lambda^2}(\lambda) \right\} (x_i/(x_i)_{p_i} - \delta_i q_i) = 0,$$

where $\tilde{\lambda}$ is a mean value between λ and λ_i , so that

$$\begin{aligned} M(\lambda_i)a_i u_1 &= \delta_i \frac{dM}{d\lambda}(\lambda_i)x_i/(x_i)_{p_i} \\ &- \delta_i^2 \left\{ \frac{dM}{d\lambda}(\lambda_i)q_i + \frac{1}{2} \frac{d^2M}{d\lambda^2}(\tilde{\lambda})(x_i/(x_i)_{p_i} - \delta_i q_i) \right\}. \end{aligned} \quad (4.2)$$

In the first stage of the iteration a vector v_{i+1} is calculated from

$$M(\lambda_i)\delta_i v_{i+1} = \delta_i \frac{dM}{d\lambda}(\lambda_i)x_i/(x_i)_{p_i} \quad (4.3)$$

so that

$$\begin{aligned} M(\lambda_i)(\delta_i v_{i+1} - a_i u_1) &= \frac{\delta_i^2}{2} \frac{d^2M}{d\lambda^2}(\lambda)x_i/(x_i)_{p_i} \\ &+ \delta_i^2 \left\{ \frac{dM}{d\lambda}(\lambda_i) - \frac{\delta_i}{2} \frac{d^2M}{d\lambda^2}(\tilde{\lambda}) \right\} q_i \\ &= \delta_i^2 y_i. \end{aligned} \quad (4.4)$$

Using equation (3.8) with $B = \frac{dM}{d\lambda}(\lambda_i) - \frac{\delta_i}{2} \frac{d^2M}{d\lambda^2}(\tilde{\lambda})$ and $\beta = T^{-1}y_i$ gives

$$\delta_i v_{i+1} - a_i u_1 = \left\{ \frac{\delta_i \beta_1}{G_{11}} + O(\delta_i^2) \right\} u_1 + r_{i+1}, \quad (4.5)$$

where $\|r_{i+1}\| = O(\delta_i^2)$. (4.6)

Thus $v_{i+1} = a_i^* u_1 + \delta_i r_{i+1}^*$ (4.7)

where $a_i^* = O(1/\delta_i)$, and $\|r_{i+1}^*\| = O(1)$.

The vector x_{i+1} is now calculated using equation (1.3). The previous argument gives

$$\begin{aligned} M(\lambda_i)(\delta_i x_{i+1} - a_i^* u_1) &= \frac{\delta_i^2}{2} \frac{d^2M}{d\lambda^2}(\lambda)v_{i+1} \\ &+ \delta_i^2 \left\{ \frac{dM}{d\lambda}(\lambda_i) - \frac{\delta_i}{2} \frac{d^2M}{d\lambda^2}(\lambda) \right\} r_{i+1}^* \\ &= \delta_i z_i. \end{aligned} \quad (4.8)$$

The right-hand side of equation (4.8) is only $O(\delta_i)$ because of the occurrence of the term involving v_{i+1} . Again applying equation (3.8) (with $\beta = T^{-1}z_i$) gives

$$\delta_i x_{i+1} - a_i^* u_1 = \left\{ \frac{\beta_1}{G_{11}} + O(\delta_i) \right\} u_1 + s_{i+1} \quad (4.9)$$

where $\|s_{i+1}\| = O(\delta_i)$. (4.10)

From equations (4.9) and (4.10) it is seen that the error in $x_{i+1}/(x_{i+1})_{p_{i+1}}$ is proportional to the square (not cube) of δ_i . Calculating the error in λ_{i+1} using equations (4.7) and (4.9) shows that this is also $O(\delta_i^2)$. Thus this iteration is only of second order.

To obtain a third-order iteration it is necessary to reduce the right-hand side of equation (4.8) to $O(\delta_i^2)$. One possibility follows from noting that, by equation (3.9),

$$a_i(u_1)_{p_i} = 1 + O(\delta_i) \quad (4.11)$$

so that, by equation (4.5),

$$(v_{i+1})_{p_i} = \frac{1}{\delta_i} + O(1). \quad (4.12)$$

If x_{i+1} is now defined to be the solution of

$$M(\lambda_i)x_{i+1} = \left[\frac{dM}{d\lambda}(\lambda_i) - \frac{1}{2(v_{i+1})_{p_i}} \frac{d^2M}{d\lambda^2}(\lambda_i) \right] v_{i+1} \quad (4.13)$$

then, using that $\lambda_i - \tilde{\lambda} = O(\delta_i)$, it is found that

$$M(\lambda_i)(\delta_i x_{i+1} - a_i^* u_1) = O(\delta_i^2) \quad (4.14)$$

so that this modification gives a third-order iteration. Equations (4.5) and (4.12) also suggest that the iteration

$$\left. \begin{aligned} M(\lambda_i)v_{i+1} &= \frac{dM}{d\lambda}(\lambda_i)v_i/(v_i)_{p_i}, \\ \lambda_{i+1} &= \lambda_i - \frac{1}{(v_{i+1})_{p_i}} \end{aligned} \right\} \quad (4.15)$$

is of second order. This iteration has the advantage of requiring the solution of only one set of linear equations per step. The author has not been able to derive this iteration by the methods used in Osborne and Michaelson.

5. Conclusion

In this paper the class of iterations discussed by Osborne and Michaelson has been examined critically with a view to determining the rate of convergence, and modifications have been suggested which ensure that the iteration is of third-order. The modified iteration (3.2) has already been tested and proved effective. The results obtained have been consistently better than those obtained using the iteration (1.3).

The iteration (4.15) has been tested on a family of

matrices all of which depended non-linearly on the eigenvalue parameter. These eigenvalue problems were also solved using the iteration (1.3). In terms of the number of iterations required there was nothing to choose between the two methods (ten runs using equation (4.15) required a total of 44 iterations, the same problems with the same starting values required a total of 43 iterations using equation (1.3)). Thus equation (4.15) provided the more efficient procedure for the problems considered.

References

- OSBORNE, M. R., and MICHAELSON, S. (1964). "The numerical solution of eigenvalue problems in which the eigenvalue parameter appears nonlinearly, with an application to differential equations," *The Computer Journal*, Vol. 7, pp. 58–65.
- WILKINSON, J. H. (1961). "Rigorous Error Bounds for Computer Eigensystems," *The Computer Journal*, Vol. 4, pp. 230–241.
- WILKINSON, J. H. (1958). "The Calculation of the Eigenvectors of Codiagonal Matrices," *The Computer Journal*, Vol. 1, pp. 90–96.

An error analysis of finite-difference methods for the numerical solution of ordinary differential equations

By M. R. Osborne*

A method is given for the calculation of strict, a-posteriori error bounds for the numerical solution by finite-difference methods of ordinary linear differential equations. The suggested procedure is illustrated by some numerical results for a particular differential equation.

1. Introduction

This paper is concerned with the derivation of strict, a-posteriori error bounds for the solution of ordinary differential equations by finite-difference methods. The basic idea of the method of error analysis is due to J. H. Wilkinson who has applied it to obtain strict error bounds for the solutions of sets of linear algebraic equations (see Wilkinson (1963) and the references quoted there). He argues as follows. Let the set of linear equations to be solved be

$$Ax = b \quad (1.1)$$

In any process of calculation rounding errors almost always occur so that the process of numerical solution leads not to x but to a vector z satisfying the system of equations

$$(A + \delta A)z = b + \delta b. \quad (1.2)$$

By a careful analysis, bounds can be found for the magnitudes of the elements of δA and δb . By suitably combining equations (1.1) and (1.2) and taking norms the result is obtained that

$$\|x - z\| \leq \frac{\|(A + \delta A)^{-1}\|}{1 - \|(A + \delta A)^{-1}\| \|\delta A\|} \{ \|\delta b\| + \|\delta A\| \|z\| \}. \quad (1.3)$$

At this stage the quantities on the right-hand side of equation (1.3) can be estimated with the exception of $\|(A + \delta A)^{-1}\|$ but, as equation (1.2) is the equation which is actually solved, there remains the possibility of obtaining at least an upper bound for the norm of the inverse of $(A + \delta A)$.

An essential feature of the argument is the representation of the vector z as the solution of a set of linear equations. This permits the original problem to be treated as a perturbation of the one actually solved once bounds have been obtained for δA and δb .

In general an error analysis of Wilkinson type does not seem to be applicable to the numerical solution of differential equations by finite-difference methods. This is because the process of solution requires the inversion of an operator of a different kind (the operator associated with a finite system of linear or non-linear algebraic equations) to the operator in the original problem.

This difficulty can be avoided in the case of finite-difference approximation to ordinary linear differential equations. Here it can be shown that the solution of the differential equation is also the solution of a linear difference equation, and bounds can be given for the difference between the coefficients in this equation and those in the equation produced by finite-difference approximation. The exact difference equation leads to

* *Edinburgh University Computer Unit, 7, Buccleuch Place, Edinburgh 8.*