

ITERATIVE METHODS FOR SOLVING MATRIX EQUATIONS

R. S. VARGA, Case Institute of Technology

1. Introduction. Iterative methods are, in concept, well known to us all since it is likely that each of us has at one time or another used Newton's method to find square roots of numbers. It is probably not as well known that iterative methods, utilizing the great speeds of modern computing machines, are extensively used today in practical computations for solving matrix equations which arise from finite difference approximations to elliptic partial differential equations of reactor technology and petroleum technology.

The object of this paper is to illustrate the use of finite difference techniques, and to illustrate the nature of iterative methods for solving the associated matrix equations. We shall do this by means of a very simple one-dimensional problem. It is hoped that this simple example will serve as an elementary introduction to the theory of such iterative methods which is covered in much more detail in [1, 2, 3, 5, 6].

2. A simple two-point boundary value problem. Consider the solution of the two-point boundary value problem:

$$(1) \quad -y^{(2)}(x) + \sigma y(x) = f(x), \quad 0 < x < 1, \quad y^{(2)} = d^2y/dx^2,$$

where

$$(2) \quad y(0) = \alpha, \quad y(1) = \beta.$$

We assume that α , β , and σ are given constants with $\sigma \geq 0$, and $f(x)$ is a given function such that $y^{(4)}(x)$ exists in $0 \leq x \leq 1$, and

$$(3) \quad |y^{(4)}(x)| \leq M, \quad 0 \leq x \leq 1.$$

By means of Taylor's Theorem, we now express $-y^{(2)}$ in terms of a three-point central difference approximation plus an error:

$$(4) \quad -y^{(2)}(x_i) = [2y(x_i) - (y(x_i + h) + y(x_i - h))]/h^2 + h^2y^{(4)}(x_i + \theta_i h)/12,$$

where $|\theta_i| < 1$, $x_i \equiv ih$, $1 \leq i \leq N$, and $h \equiv 1/(N+1)$. With $y(x_i) \equiv y_i$, the differential equation (1) can be written for the particular values x_i , $1 \leq i \leq N$, in matrix form, as

$$(5) \quad Ay = k + \tau(y),$$

where A is an $N \times N$ real matrix, and \mathbf{y} , \mathbf{k} , and $\boldsymbol{\tau}(\mathbf{y})$ are vectors, all given explicitly by

$$(6) \quad A = \frac{1}{h^2} \begin{bmatrix} 2 + \sigma h^2 & & -1 & & \\ & -1 & & 2 + \sigma h^2 & \\ & & -1 & & -1 \\ & & & -1 & & 2 + \sigma h^2 \\ & & & & -1 & & 2 + \sigma h^2 \end{bmatrix}; \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}; \quad \mathbf{k} = \begin{bmatrix} f_1 + \alpha/h^2 \\ f_2 \\ \vdots \\ f_N + \beta/h^2 \end{bmatrix},$$

$$\boldsymbol{\tau} = \frac{-h^2}{12} \begin{bmatrix} y^{(4)}(x_1 + \theta_1 h) \\ \vdots \\ y^{(4)}(x_N + \theta_N h) \end{bmatrix}.$$

Neglecting the vector $\boldsymbol{\tau}(\mathbf{y})$ in (5) gives us the matrix equation

$$(7) \quad A\mathbf{z} = \mathbf{k},$$

whose solution \mathbf{z} is defined to be our discrete approximation to the solution $y(x)$ of (1)–(2), i.e., z_j , the j th component of \mathbf{z} , is to approximate $y_j = y(jh)$.

It is obvious that A is real and symmetric. Noting that all the diagonal entries of A are equal, we can express the matrix A as

$$(8) \quad A = \frac{(2 + \sigma h^2)}{h^2} [I - B],$$

where B is an $N \times N$ real symmetric matrix given explicitly by

$$B = \frac{1}{(2 + \sigma h^2)} \begin{bmatrix} 0 & 1 & & 0 \\ 1 & 0 & 1 & \\ & 1 & 0 & 1 \\ & & 1 & 0 & 1 \\ 0 & & & 1 & 0 \end{bmatrix}.$$

Because of the tridiagonal form of B , it is easy to verify that the vector $\mathbf{x}^{(j)}$, with k -th component $x_k^{(j)}$ defined by $x_k^{(j)} = \sin(jk\pi h)$, $1 \leq k \leq N$, is an eigenvector of B for each j , with the corresponding eigenvalue μ_j given by

$$(10) \quad \mu_j = 2 \cos(j\pi h)/(2 + \sigma h^2), \quad 1 \leq j \leq N.$$

Then, it follows that the spectral radius $\rho(B) \equiv \max_{1 \leq j \leq N} |\mu_j|$ is less than unity, and thus, from (8), A is evidently positive definite. Moreover, since $\rho(B) < 1$, we can express A^{-1} by means of the convergent matrix infinite series

$$(11) \quad A^{-1} = \gamma(I - B)^{-1} = \gamma\{I + B + B^2 + \dots\}, \quad \gamma \equiv h^2/(2 + \sigma h^2).$$

Note that all entries of B are nonnegative from (9), so that all powers of B and hence A^{-1} have only nonnegative real entries. This last fact will be useful in proving in the next section that the solution \mathbf{z} of (7) is "close" to the solution $y(x)$ of (1)–(2), evaluated at the points $x_i = ih$, $1 \leq i \leq N$.

3. Convergence of the discrete approximation. To measure the discrepancy between the components z_i of the vector \mathbf{z} of (7) and the numbers $y_i = y(ih)$, we subtract (7) from (5) giving $A(\mathbf{y} - \mathbf{z}) = \boldsymbol{\tau}(\mathbf{y})$, or, since A is nonsingular,

$$(12) \quad \mathbf{y} - \mathbf{z} = A^{-1}\boldsymbol{\tau}(\mathbf{y}).$$

Now, from (3), each component of $\boldsymbol{\tau}(\mathbf{y})$ is bounded above by $Mh^2/12$. Hence, since A^{-1} has only nonnegative components, $|y_i - z_i| \leq (Mh^2/12)(A^{-1}\boldsymbol{\xi})_i$ for all $1 \leq i \leq N$, where $\boldsymbol{\xi}$ is the vector with all components unity. Thus, if \mathbf{w} is any vector with $A\mathbf{w} \geq \boldsymbol{\xi}$, we obtain the bound $|y_i - z_i| \leq Mh^2w_i/12$. But it is easy to verify that the choice $w(x) \equiv x(1-x)/2$ with $w_i = w(ih)$ satisfies $A\mathbf{w} \geq \boldsymbol{\xi}$, so that

$$(13) \quad |y_i - z_i| \leq Mh^2[ih(1 - ih)]/24 \leq Mh^2/96,$$

and last inequality following since $w(x) \leq 1/8$ for $0 \leq x \leq 1$. This means that \mathbf{z} is close to \mathbf{y} if the mesh spacing h is sufficiently small.

4. The Jacobi iterative method. For h small, the result of the previous section shows that solving $A\mathbf{z} = \mathbf{k}$ will give us a reasonable approximation to the solution of the original differential equation (1)-(2). To solve $A\mathbf{z} = \mathbf{k}$ for \mathbf{z} , let us combine (7) with (8) to form

$$(14) \quad \mathbf{z} = B\mathbf{z} + \mathbf{g},$$

where $\mathbf{g} \equiv (h^2\mathbf{k})/(2 + \sigma h^2)$. Familiarity with the method of successive substitutions suggests that we consider the following iterative method, called the Jacobi iterative method,

$$(15) \quad \mathbf{z}^{(m+1)} = B\mathbf{z}^{(m)} + \mathbf{g}, \quad m \geq 0,$$

where $\mathbf{z}^{(0)}$ is some initial estimate of the unique solution of (14). Writing $\mathbf{z}^{(m)} = \mathbf{z} + \boldsymbol{\epsilon}^{(m)}$ to define the error $\boldsymbol{\epsilon}^{(m)}$ at each iteration, we see from (14) and (15) that $\boldsymbol{\epsilon}^{(m+1)} = B\boldsymbol{\epsilon}^{(m)}$, from which we deduce that

$$(16) \quad \boldsymbol{\epsilon}^{(m)} = B^m\boldsymbol{\epsilon}^{(0)}, \quad m \geq 0.$$

Since B is a real symmetric $N \times N$ matrix, the eigenvectors $\{\mathbf{x}^{(j)}\}_{j=1}^N$ of B can be normalized to form an orthonormal basis for the associated N dimensional vector space. Thus, there exist constants c_j such that $\boldsymbol{\epsilon}^{(0)} = \sum_{j=1}^N c_j \mathbf{x}^{(j)}$, where $B\mathbf{x}^{(j)} = \mu_j \mathbf{x}^{(j)}$, and it follows from (16) that

$$(17) \quad \boldsymbol{\epsilon}^{(m)} = \sum_{j=1}^N (\mu_j)^m c_j \mathbf{x}^{(j)}, \quad m \geq 0.$$

Because all μ_j are less than unity in modulus, $\boldsymbol{\epsilon}^{(m)}$ evidently tends to the zero vector for any initial $\boldsymbol{\epsilon}^{(0)}$. Equivalently, the Jacobi iterative method converges for any initial $\mathbf{z}^{(0)}$. What is more, from (10) each component of the error $\boldsymbol{\epsilon}^{(m)}$ with respect to the orthonormal basis $\{\mathbf{x}^{(j)}\}_{j=1}^N$ is reduced per iteration by a factor not exceeding in modulus

$$(18) \quad \rho(B) = 2 \cos \pi h / (2 + \sigma h^2) = 1 - (\pi^2 + \sigma)h^2/2 + O(h^4), \quad h \downarrow 0.$$

Having just shown that the iterative method of (15) is convergent, it is our sad duty to report that this iterative method can be *extremely* slowly convergent for small h . To illustrate this, let $h = 10^{-2}$ and $\sigma = 1$. Then, to reduce each component of $\mathbf{e}^{(0)}$ by 10 would require approximately m iterations, where $[\rho(B)]^m \approx 0.1$, and m turns out to be about 4200. Even for fast computing machines, this would be a slow process, and obviously there is a need for faster iterative methods.

5. The successive overrelaxation iterative method. To introduce just one such faster iterative method, we first express the matrix B of (9) as the sum $L + L^T$, where L is a strictly lower triangular matrix. Multiplying both sides of (14) by the real parameter ω , and then adding \mathbf{z} to both sides of the result yields, after rearrangement,

$$(19) \quad (I - \omega L)\mathbf{z} = \{(1 - \omega)I + \omega L^T\}\mathbf{z} + \omega \mathbf{g}.$$

Since L is strictly lower triangular, then $(I - \omega L)$ is nonsingular for any choice of ω . Thus, we can multiply both sides of (19) on the left by $(I - \omega L)^{-1}$, which gives

$$(19') \quad \mathbf{z} = (I - \omega L)^{-1}[(1 - \omega)I + \omega L^T]\mathbf{z} + \omega(I - \omega L)^{-1}\mathbf{g}.$$

Simply inserting iteration superscripts then serves to define the successive overrelaxation (SOR) iterative method:

$$(20) \quad \mathbf{z}^{(m+1)} = (I - \omega L)^{-1}[(1 - \omega)I + \omega L^T]\mathbf{z}^{(m)} + \omega(I - \omega L)^{-1}\mathbf{g}, \quad m \geq 0.$$

At first glance, this iterative method of (20) looks rather formidable, and would appear to be implicit. From the definition of B in (9), however, this iterative method can also be written equivalently as

$$(21) \quad z_i^{(m+1)} = z_i^{(m)} - \omega[(2 + \sigma h^2)z_i^{(m)} - z_{i-1}^{(m+1)} - z_{i+1}^{(m)} - h^2 f_i]/(2 + \sigma h^2),$$

$$1 \leq i \leq N.$$

Starting with $i=1$, $z_0 = \alpha$ is fixed by the boundary condition (2), and we can solve for the new component $z_1^{(m+1)}$ in (21) in terms of the old components $z_1^{(m)}$ and $z_2^{(m)}$. Continuing from left to right, we see that this is an *explicit* iterative method.

In analogy with the Jacobi iterative method, we similarly seek the eigenvalues λ of the matrix

$$(22) \quad \mathcal{L}_\omega \equiv (I - \omega L)^{-1}\{(1 - \omega)I + \omega L^T\}.$$

Thus, we seek the roots of $\det\{\lambda I - \mathcal{L}_\omega\} = 0$. Since $\det(I - \omega L) = 1$, it follows that

$$(23) \quad 0 = \det(I - \omega L) \det(\lambda I - \mathcal{L}_\omega) = \det[(\lambda + \omega - 1)I - \omega(\lambda L + L^T)].$$

Previously, we explicitly determined the eigenvectors $\mathbf{x}^{(j)}$ and eigenvalues μ_j of the matrix B of (9). In the same manner, one can verify that the vector $\mathbf{w}^{(j)}$, with k th component $w_k^{(j)}$ defined by $\tau^{k/2} \sin(jk\pi h)$, is an eigenvector of $\tau L + L^T$ for any τ , with corresponding eigenvalue $\tau^{1/2}\mu_j$. Consequently, with $\tau = \lambda$ there is a natural pairing from (23) of the eigenvalues λ of \mathcal{L}_ω with the eigenvalues μ of B through

$$(24) \quad (\lambda + \omega - 1)^2 = \lambda\omega^2\mu^2,$$

which is the fundamental result of Young [8]. The real significance of this expression (24) is that it can be used to explicitly determine the real parameter ω which *minimizes* the spectral radius of \mathcal{L}_ω . Indeed, it is now well known [8] that

$$(25) \quad \min_{\omega} \rho(\mathcal{L}_\omega) = \rho(\mathcal{L}_{\omega_b}) = \omega_b - 1, \quad \text{where } \omega_b = 2/(1 + \sqrt{1 - \rho^2(B)}).$$

Using the result of (18), we see that

$$(26) \quad \rho(\mathcal{L}_{\omega_b}) = 1 - 2(\sqrt{\pi^2 + \sigma})h + O(h^2), \quad h \downarrow 0.$$

Comparison of the exponents of h in the leading terms of (18) and (26) shows that an *order of magnitude* improvement in convergence rate has been made, simply by the adroit choice of the parameter ω . To illustrate this numerically, again let $h = 10^{-2}$ and $\sigma = 1$. Then, the least positive integer m such that $[\rho(\mathcal{L}_{\omega_b})]^m \leq 0.1$ now turns out to be approximately 35, as compared with 4200. For the slight increase in arithmetic requirements per iteration in passing from (15) to (21), a great over-all improvement in computational efficiency has been achieved.

6. Alternating direction implicit methods. After all is said and done, the solution of the discrete one-dimensional problem $Az = \mathbf{k}$ in (7) would *not* be determined by means of iterative methods, except, of course, in expository papers on the subject! Rather, straight-forward Gaussian elimination applied to the matrix problem $Az = \mathbf{k}$ is not only efficient, but it is very stable relative to the growth of rounding errors, thanks to the positive definite tridiagonal nature of A . The Gaussian elimination algorithm for solving this simple problem can be expressed as

$$(27) \quad \begin{cases} w_1 = -1/b; & w_i = (-1/(b + w_{i-1})), \quad 2 \leq i \leq N-1, \quad b \equiv 2 + \sigma h^2, \\ g_1 = h^2 k_1/b; & g_i = (h^2 k_i + g_{i-1})/(b + w_{i-1}), \quad 2 \leq i \leq N, \\ z_N = g_N; & z_i = g_i - w_i z_{i+1}, \quad 1 \leq i \leq N-1. \end{cases}$$

Our simple one-dimensional problem can be useful once more to us in quickly introducing a particular variant of the alternating direction implicit (ADI) methods, for the discussion above shows that discrete approximations to the one-dimensional problem $-u_{xx} + \sigma u = f$ can be directly solved. Consider then

the second order elliptic partial differential equation

$$(28) \quad -u_{xx}(x, y) - u_{yy}(x, y) + 2\sigma u(x, y) = f(x, y), \quad 0 < x, y < 1,$$

in the unit square R with Dirichlet boundary conditions

$$(29) \quad u(x, y) = g(x, y), \quad (x, y) \in \partial R,$$

where g is specified on ∂R , the boundary of R . Writing (28) as

$$(30) \quad [-u_{xx} + \sigma u] + [-u_{yy} + \sigma u] = f,$$

each term in brackets represents a differential operator in one of the space variables. Thus, with

$$(31) \quad \begin{cases} Hu(x_0, y_0) \equiv [(2 + \sigma h^2)u(x_0, y_0) - (u(x_0 + h, y_0) + u(x_0 - h, y_0))]/h^2, \\ Vu(x_0, y_0) \equiv [(2 + \sigma h^2)u(x_0, y_0) - (u(x_0, y_0 + h) + u(x_0, y_0 - h))]/h^2, \end{cases}$$

representing discrete approximations to these differential operators on a uniform mesh $h = 1/(N+1)$, the discrete matrix problem corresponding to (28)–(29) is

$$(32) \quad Hu + Vu = \mathbf{k},$$

where H and V are real $N^2 \times N^2$ matrices. It is not difficult to see that, with suitable numbering of the mesh points of the square, H and V are direct sums of the particular matrix A of (6). An important application of the discussion above concerning Gaussian elimination is that matrix problems of the form $H\mathbf{x} = \mathbf{s}$ and $V\mathbf{x} = \mathbf{t}$ can be solved *directly* for \mathbf{x} , which is to be a basic part of the iterative method to be described. Writing (32) as the pair of equations

$$(33) \quad (H + rI)\mathbf{u} = (rI - V)\mathbf{u} + \mathbf{k}, \quad (V + rI)\mathbf{u} = (rI - H)\mathbf{u} + \mathbf{k},$$

we insert iteration superscripts to define the Peaseman-Rachford alternating direction iterative method [4]

$$(34) \quad \begin{aligned} (H + r_m I)\mathbf{u}^{(m+1/2)} &= (r_m I - V)\mathbf{u}^{(m)} + \mathbf{k}, \\ (V + r_m I)\mathbf{u}^{(m+1)} &= (r_m I - H)\mathbf{u}^{(m+1/2)} + \mathbf{k}, \quad m \geq 0, \end{aligned}$$

where $\{r_m\}$ is any sequence of positive acceleration parameters. Note that carrying out a single complete iteration of (34) requires the direct solution of matrix equations on horizontal mesh lines, then on vertical mesh lines, hence the name *alternating directions*.

For the particular problem (28)–(29), we shall now show that the iterative method (34) converges for *any* fixed positive acceleration parameter $r > 0$. In analogy with section 4, we first exhibit explicitly the eigenvectors of the $N^2 \times N^2$ matrices H and V . Defining the column vector $\alpha^{(k,l)}$ with N^2 components $\alpha_{i,j}^{(k,l)}$ by

$$(35) \quad \alpha_{i,j}^{(k,l)} = \sin(k\pi ih) \sin(l\pi jh), \quad 1 \leq i, j \leq N, \quad 1 \leq k, l \leq N,$$

where $\alpha_{i,j}^{(k,l)}$ refers to the component of $\alpha^{(k,l)}$ at the i th column and j th row of the uniform mesh on the unit square, we see from (31) that

$$(36) \quad \begin{aligned} H\alpha^{(k,l)} &= [4 \sin^2(k\pi h/2) + \sigma h^2]\alpha^{(k,l)}; \\ V\alpha^{(k,l)} &= [4 \sin^2(l\pi h/2) + \sigma h^2]\alpha^{(k,l)} \quad 1 \leq k, l \leq N. \end{aligned}$$

For any fixed $r > 0$, the error vectors $\epsilon^{(m)}$ associated with (34) evidently satisfy, in analogy with (16),

$$(37) \quad \epsilon^{(m+1)} = T_r \epsilon^{(m)}, \quad m \geq 0, \quad \text{where } T_r \equiv (V + rI)^{-1}(rI - H)(rI + H)^{-1}(rI - V).$$

Thus, from (36), we have that $T_r \alpha^{(k,l)} = \lambda_{k,l} \alpha^{(k,l)}$ where

$$\lambda_{k,l} = \frac{[r - (4 \sin^2(k\pi h/2) + \sigma h^2)][r - (4 \sin^2(l\pi h/2) + \sigma h^2)]}{[r + (4 \sin^2(k\pi h/2) + \sigma h^2)][r + (4 \sin^2(l\pi h/2) + \sigma h^2)]}, \quad 1 \leq k, l \leq N,$$

and clearly, $|\lambda_{k,l}| < 1$ for any $1 \leq k, l \leq N$ since $r > 0$. Because the $\alpha^{(k,l)}$ form, after normalization, an orthonormal basis for the associated N^2 dimensional vector space, it follows from $|\lambda_{k,l}| < 1$ that $\epsilon^{(m)}$ tends to the zero vector as $m \rightarrow \infty$ for any initial vector $\epsilon^{(0)}$, which proves that the iterative method of (34) converges for any fixed $r > 0$. If we now choose the parameter r to be

$$r = \{(4 \sin^2(\pi h/2) + \sigma h^2)(4 \cos^2(\pi h/2) + \sigma h^2)\}^{1/2},$$

it turns out [6, p. 216] that the spectral radius $\rho(T_r)$ is exactly equal to that of the successive overrelaxation iterative method applied to (28)–(29), so that from (26),

$$(38) \quad \rho(T_r) = 1 - 2(\sqrt{\pi^2 + \sigma})h + O(h^2), \quad h \downarrow 0.$$

The really interesting results for such ADI methods are obtained by using a sequence of optimally selected positive parameters [7; 6, p. 223], since the use of such parameters can lead to a further order of magnitude improvement in convergence rates. Specifically, with such optimally selected parameters, we merely state that the average error reduction in norm per iteration turns out to be asymptotically

$$(39) \quad 1 - \frac{(\pi^2 + \sigma)}{\ln(1/h)}, \quad h \downarrow 0,$$

which should be contrasted with (38). It must be mentioned, however, that the result of (39) depends heavily on the assumption that the matrices H and V have the common set of eigenvectors $\alpha^{(k,l)}$, which is equivalent to the assumption that H and V commute, i.e., $HV = VH$. For cases in which $HV \neq VH$, the results concerning convergence rates of ADI methods are much less complete.

There is a rich and growing literature on the subject of iterative methods; our aim here was to introduce quickly some of these methods. More complete results, beyond the scope of this paper, can be found in [2, 3, 5, 6].

References

1. D. K. Faddeev and V. N. Faddeeva, Computational methods of linear algebra, translated by R. C. Williams, Freeman, San Francisco, 1963.
2. G. E. Forsythe and W. R. Wasow, Finite difference methods for partial differential equations, Wiley, New York, 1960.
3. A. S. Householder, The theory of matrices in numerical analysis, Blaisdell, New York, 1964.
4. D. W. Peaceman and H. H. Rachford, Jr., The numerical solution of parabolic and elliptic differential equations, *J. Soc. Indust. Appl. Math.*, 3 (1955) 28-41.
5. John Todd, (editor) Survey of numerical analysis, McGraw-Hill, New York, 1962.
6. R. S. Varga, Matrix iterative analysis, Prentice-Hall, Englewood Cliffs, N. J., 1962.
7. E. L. Wachspress, Optimum alternating-direction-implicit iteration parameters for a model problem, *J. Soc. Indust. Appl. Math.*, 10 (1962) 339-350.
8. D. M. Young, Iterative methods for solving partial difference equations of elliptic type, *Trans. Amer. Math. Soc.*, 76 (1954) 92-111.