

7 Iterative methods for matrix equations

7.1 The need for iterative methods

We have seen that Gaussian elimination provides a method for finding the exact solution (if rounding errors can be avoided) of a system of equations $\mathbf{Ax} = \mathbf{b}$. However Gaussian elimination requires approximately $n^3/3$ operations (where n is the size of the system), which may become prohibitively time-consuming if n is very large. Another weakness is that Gaussian elimination requires us to store all the components of the matrix \mathbf{A} . In many real applications (especially the numerical solution of differential equations), the matrix \mathbf{A} is *sparse*, meaning that most of its elements are zero, in which case keeping track of the whole matrix is wasteful.

In situations like these it may be preferable to adopt a method which produces an approximate rather than exact solution. We will describe three *iterative methods*, which start from an initial guess \mathbf{x}_0 and produce successively better approximations $\mathbf{x}_1, \mathbf{x}_2, \dots$. The iteration can be halted as soon as an adequate degree of accuracy is obtained, and the hope is that this takes a significantly shorter time than the exact method of Gaussian elimination would require.

7.2 Splitting the matrix

All the methods we will consider involve splitting the matrix \mathbf{A} into the difference between two new matrices \mathbf{S} and \mathbf{T} :

$$\mathbf{A} = \mathbf{S} - \mathbf{T}.$$

Thus the equation $\mathbf{Ax} = \mathbf{b}$ gives

$$\mathbf{Sx} = \mathbf{Tx} + \mathbf{b},$$

based on which we can try the iteration

$$\mathbf{Sx}_{k+1} = \mathbf{Tx}_k + \mathbf{b}. \tag{7.1}$$

Now *if* this procedure converges, say $\mathbf{x}_k \rightarrow \mathbf{x}$ as $k \rightarrow \infty$, then clearly \mathbf{x} solves the original problem $\mathbf{Ax} = \mathbf{b}$, but it is not at all clear from the outset whether a scheme like (7.1) converges or not.

Evidently there are many possible ways to split the matrix \mathbf{A} . The tests of a good choice are:

- the new vector \mathbf{x}_{k+1} should be *easy to compute*, that is \mathbf{S} should be easily invertible (for example \mathbf{S} might be diagonal or triangular);
- the scheme should *converge as rapidly as possible* towards the true solution.

These two requirements are conflicting: a choice of splitting which is particularly easy to invert (see *e.g.* Jacobi's method below) may not converge especially rapidly (or at all). At the other extreme we can converge *exactly*, in just one step, by using $\mathbf{S} = \mathbf{A}$, $\mathbf{T} = \mathbf{0}$; but $\mathbf{S} = \mathbf{A}$ is usually difficult to invert: that's the whole point of splitting!

It is convenient to introduce the notation

$$\mathbf{A} = \mathbf{L} + \mathbf{D} + \mathbf{U} \quad (= \mathbf{S} - \mathbf{T}),$$

where \mathbf{L} is strictly lower triangular,
 \mathbf{D} is diagonal,
 \mathbf{U} is strictly upper triangular.

For example, if

$$\mathbf{A} = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix},$$

then

$$\mathbf{L} = \begin{pmatrix} 0 & 0 & 0 \\ 4 & 0 & 0 \\ 7 & 8 & 0 \end{pmatrix}, \quad \mathbf{D} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 9 \end{pmatrix}, \quad \mathbf{U} = \begin{pmatrix} 0 & 2 & 3 \\ 0 & 0 & 6 \\ 0 & 0 & 0 \end{pmatrix}.$$

We will consider three possible splittings.

1. **Jacobi's method** $\mathbf{S} = \mathbf{D}$, $\mathbf{T} = -(\mathbf{L} + \mathbf{U})$;
2. **The Gauss-Seidel method** $\mathbf{S} = \mathbf{L} + \mathbf{D}$, $\mathbf{T} = -\mathbf{U}$;
3. **Successive over-relaxation (SOR)** a combination of 1 and 2.

7.3 Jacobi's method

In Jacobi's method, \mathbf{S} is simply the diagonal part of \mathbf{A} . We illustrate it with a simple two-dimensional example.

Example 7.1 Consider the system

$$\begin{aligned} 2x - y &= 3, \\ -x + 2y &= 0. \end{aligned}$$

We solve the first equation for x and the second for y :

$$\begin{aligned} x &= y/2 + 3/2, \\ y &= x/2, \end{aligned}$$

which suggests the iterative scheme

$$\begin{aligned}x_{k+1} &= y_k/2 + 3/2 \\y_{k+1} &= x_k/2.\end{aligned}$$

Note that this corresponds to $\mathbf{S}\mathbf{x}_{k+1} = \mathbf{T}\mathbf{x}_k + \mathbf{b}$ where the splitting

$$\mathbf{S} = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}, \quad \mathbf{T} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix},$$

has been employed.

Suppose we start with the initial guess

$$\begin{pmatrix} x_0 \\ y_0 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}.$$

Then we generate successively

$$\begin{pmatrix} x_1 \\ y_1 \end{pmatrix} = \begin{pmatrix} 2 \\ 1/2 \end{pmatrix}, \quad \begin{pmatrix} x_2 \\ y_2 \end{pmatrix} = \begin{pmatrix} 7/4 \\ 1 \end{pmatrix}, \quad \begin{pmatrix} x_3 \\ y_3 \end{pmatrix} = \begin{pmatrix} 2 \\ 7/8 \end{pmatrix},$$

and so forth, which appears to be converging nicely to the solution $x = 2, y = 1$.

7.4 The Gauss-Seidel method

A drawback with Jacobi's method is that it requires us to store all the components of \mathbf{x}_k until we have finished computing the next iteration \mathbf{x}_{k+1} . For very large systems, the memory required may become a problem. In any case, it would appear to make more sense to use the "new" values of \mathbf{x} (which presumably are more accurate than the old values) as soon as we've calculated them. The scheme which results from doing so is called the Gauss-Seidel method. We illustrate it with the same two-dimensional system as in example 7.1.

Example 7.2 Consider the same system

$$\begin{aligned}2x - y &= 3 \\ -x + 2y &= 0,\end{aligned}$$

as in example 7.1. As in Jacobi's method we use the first equation to find x_{k+1} in terms of y_k :

$$x_{k+1} = y_k/2 + 3/2.$$

But now that we've found x_{k+1} we use it when working out y_{k+1} :

$$y_{k+1} = x_{k+1}/2$$

(Jacobi has $x_k/2$ on the right-hand side). In terms of splitting this corresponds to the choice

$$\mathbf{S} = \begin{pmatrix} 2 & 0 \\ -1 & 2 \end{pmatrix}, \quad \mathbf{T} = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix},$$

This form persists for the subsequent steps, so the general formula for updating x_i is

$$(x_i)_{\text{new}} = x_i + \frac{1}{a_{ii}} \left[b_i - \sum_{j=1}^n a_{ij}x_j \right]. \quad (7.2)$$

On the right-hand side of equation (7.2), the first term is the “old” value of x_i , and the second term is the *correction*, which is added to x_i to obtain the improved value $(x_i)_{\text{new}}$. Note also from the form of (7.2) that for the scheme to work *we must set up the system in such a way that its diagonal entries are all nonzero*.

We can estimate the number of operation required for the Gauss-Seidel method as follows. Each step of the form (7.2) requires n multiplications and 1 division, making a total of $n + 1$ operations. These are required for each of the n components of \mathbf{x} , so that $n(n + 1) \approx n^2$ operations are needed for one complete iteration. Thus the total number of operations required is approximately n^2 times the number of iterations required for satisfactory convergence. This should be compared with the $n^3/3$ operations needed for Gaussian elimination. We can observe that Gauss-Seidel (or something similar) may indeed be preferable to Gaussian elimination if n is very large and *the number of iterations required can be kept reasonably low*. We should also observe that in the operations count carried out above, the number of multiplications required may be greatly reduced if the matrix A is sparse.

The question arises: *how can we tell if the scheme has converged adequately?* That is, how do we decide when enough iterations have been performed? One possibility is to look at the difference between successive values of \mathbf{x}_k , since these become closer as \mathbf{x}_k approaches the solution \mathbf{x} . Thus we might, for example, choose to continue the iterations until $|\mathbf{x}_{k+1} - \mathbf{x}_k|$ is smaller than some prescribed small amount. A slightly better approach is to use the *relative* difference (*i.e.* compared to the actual size of the solution) as a measure of how well the scheme has converged. Thus the iteration is halted when

$$\frac{|\mathbf{x}_{k+1} - \mathbf{x}_k|}{|\mathbf{x}_{k+1}|} < \delta,$$

where δ is some prescribed small number. In the Gauss-Seidel step this is very easy to calculate, since $(x_i)_{k+1} - (x_i)_k$ is just the second term on the right-hand side of (7.2).

We implement these ideas in the following fragment of code in FORTRAN90 which performs a single iteration of the Gauss-Seidel scheme.

```
!
!   This code performs a single Gauss-Seidel iteration
!
error = 0.0
do i = 1, n
  temp = b(i)
  do j = 1, n
    temp = temp - a(i,j) * x(j)
  end do
```

```

        temp = temp / a(i,i)
!
!       Now temp contains the correction to be added to x(i),
!       that is (x(i))_(k+1) - (x(i))_k
!
        error = error + temp**2
        x(i) = x(i) + temp
    end do
    error = sqrt(error)
!
!       Now error contains |x_(k+1) - x_k|
!

```

All that remains is to perform this loop until the relative error is sufficiently small. We should also allow for the fact that the scheme may not converge at all. The following code demonstrates a straightforward implementation, where the “iteration step” is the code fragment given above.

```

... preamble ...
    do k = 1, 30
... iteration step ...
        if (error / sqrt(DOT_PRODUCT(x,x)) < 0.001) exit
    end do
    if (k > 30) then                ! Scheme hasn't converged
        write(*,*) 'Method failed to converge after 30 iterations'
    else                            ! Scheme has converged
        write(*,*) 'Method converged after', k, 'iterations'
        write(*,*) 'Solution is'
        do i = 1, n
            write(*,*) x(i)
        end do
    end if
... etc. ...

```

Note that here we use the function DOT_PRODUCT to calculate $|\mathbf{x}_k|$.

7.6 Convergence of Jacobi and Gauss-Seidel

In general our iterative process takes the form (7.1). Of course, for the method to be at all sensible \mathbf{S} must be invertible, so we can write

$$\mathbf{x}_{k+1} = \mathbf{S}^{-1}\mathbf{T}\mathbf{x}_k + \mathbf{S}^{-1}\mathbf{b}.$$

The combination $(\mathbf{S}^{-1}\mathbf{T})$ is known as the *iteration matrix*.

Now we set

$$\mathbf{x}_k = \mathbf{x} + \boldsymbol{\epsilon}_k,$$

where \mathbf{x} is the exact solution (*i.e.* it satisfies $\mathbf{x} = \mathbf{S}^{-1}\mathbf{T}\mathbf{x} + \mathbf{S}^{-1}\mathbf{b}$) and $\boldsymbol{\epsilon}_k$ is the error after k iterations. Then $\boldsymbol{\epsilon}_k$ satisfies

$$\boldsymbol{\epsilon}_{k+1} = (\mathbf{S}^{-1}\mathbf{T}) \boldsymbol{\epsilon}_k,$$

and the question is *does* $\boldsymbol{\epsilon}_k \rightarrow \mathbf{0}$ as $k \rightarrow \infty$?

We will answer this question under the assumption that *the iteration matrix* $\mathbf{S}^{-1}\mathbf{T}$ has a complete set of n linearly independent eigenvectors, say $\mathbf{e}_1, \dots, \mathbf{e}_n$, with corresponding eigenvalues $\lambda_1, \dots, \lambda_n$. Thus we can use $\{\mathbf{e}_i\}$ as a basis and write

$$\boldsymbol{\epsilon}_0 = \sum_{i=1}^n c_i \mathbf{e}_i,$$

for some constants c_i , where $\boldsymbol{\epsilon}_0$ is the error in the initial guess \mathbf{x}_0 . Then

$$\boldsymbol{\epsilon}_1 = \mathbf{S}^{-1}\mathbf{T}\boldsymbol{\epsilon}_0 = \sum_{i=1}^n c_i \lambda_i \mathbf{e}_i,$$

and in general

$$\boldsymbol{\epsilon}_k = \sum_{i=1}^n c_i \lambda_i^k \mathbf{e}_i. \quad (7.3)$$

So $\boldsymbol{\epsilon}_k$ tends to zero as $k \rightarrow \infty$ if and only if $|\lambda_i| < 1$ for all i . That is, the iteration (7.1) converges if and only if *the eigenvalues of the iteration matrix* $\mathbf{S}^{-1}\mathbf{T}$ all have absolute value less than one.

Another way of stating this result is as follows. Define the *spectral radius* $\rho(\mathbf{M})$ of a matrix \mathbf{M} to be the largest value of $|\lambda_i|$ over the eigenvalues λ_i of \mathbf{M} :

$$\rho(\mathbf{M}) = \max_i |\lambda_i|.$$

Then the iterative scheme (7.1) converges if and only if $\rho(\mathbf{S}^{-1}\mathbf{T}) < 1$.

Another important aspect of the spectral radius is its implications for the rate at which the scheme converges. If (without loss of generality) λ_1 is the largest eigenvalue of $\mathbf{S}^{-1}\mathbf{T}$ in absolute value (with $|\lambda_1| = \rho$) then clearly (7.3) implies that the behaviour of the error $\boldsymbol{\epsilon}_k$ is dominated by λ_1^k as $k \rightarrow \infty$:

$$\begin{aligned} \boldsymbol{\epsilon}_k &\sim c_1 \lambda_1^k \mathbf{e}_1 \text{ as } k \rightarrow \infty \\ \Rightarrow |\boldsymbol{\epsilon}_k| &\sim \text{const. } \rho^k \text{ as } k \rightarrow \infty. \end{aligned}$$

Thus *the error is reduced by roughly a factor of ρ at each iteration* and so *the smaller ρ is, the faster the scheme converges*.

With this in mind, the *rate of convergence* is often defined as $-\log \rho$ when $\rho < 1$.

Finally we reiterate that we have assumed throughout this section that the iteration matrix has a complete set of eigenvectors. This may seem at first glance to be a dangerous assumption, but the following points should be borne in mind.

- A matrix is guaranteed to have a complete set of eigenvectors so long as all its eigenvalues are distinct.
- Thus the case in which two or eigenvalues are repeated, in which there is the possibility that the set of eigenvectors might be incomplete, is not generic.
- In any case, the results of this section *can*, with a little more work, be shown to hold even when $\mathbf{S}^{-1}\mathbf{T}$ does not have a complete set of eigenvectors.

Now we illustrate the calculation of the spectral radius for the Jacobi and Gauss-Seidel schemes considered in examples 7.1 and 7.2.

Example 7.3 *The Jacobi iteration from example 7.1,*

$$\begin{aligned}x_{k+1} &= y_k/2 + 3/2, \\y_{k+1} &= x_k/2,\end{aligned}$$

can be written as $\mathbf{x}_{k+1} = \mathbf{S}^{-1}\mathbf{T}\mathbf{x}_k + \mathbf{S}^{-1}\mathbf{b}$, where

$$\mathbf{S}^{-1}\mathbf{T} = \begin{pmatrix} 0 & 1/2 \\ 1/2 & 0 \end{pmatrix}, \quad \mathbf{S}^{-1}\mathbf{b} = \begin{pmatrix} 3/2 \\ 0 \end{pmatrix}.$$

The eigenvalues λ of $\mathbf{S}^{-1}\mathbf{T}$ are given by

$$\begin{vmatrix} -\lambda & 1/2 \\ 1/2 & -\lambda \end{vmatrix} = 0 \Rightarrow \lambda^2 = 1/4 \Rightarrow \lambda = \pm 1/2.$$

Thus the spectral radius $\rho_{\text{Jacobi}} = 1/2 < 1$ and the scheme therefore converges.

Example 7.4 *For the Gauss-Seidel iteration from example 7.2, we have*

$$\mathbf{S} = \begin{pmatrix} 2 & 0 \\ -1 & 2 \end{pmatrix}, \quad \mathbf{T} = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} \Rightarrow \mathbf{S}^{-1}\mathbf{T} = \frac{1}{4} \begin{pmatrix} 2 & 0 \\ 1 & 2 \end{pmatrix} \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} 0 & 1/2 \\ 0 & 1/4 \end{pmatrix},$$

using the usual formula for the inverse of a 2×2 matrix. The eigenvalues are given by

$$\begin{vmatrix} -\lambda & 1/2 \\ 0 & 1/4 - \lambda \end{vmatrix} = 0 \Rightarrow \lambda(\lambda - 1/4) = 0 \Rightarrow \lambda = 0 \text{ or } 1/4,$$

and the spectral radius $\rho_{\text{G-S}} = 1/4$.

Thus for the system considered in these examples the Gauss-Seidel scheme converges *twice as fast* as the Jacobi scheme. Recall that the error is roughly speaking reduced by a factor of ρ at each iteration. Thus it takes two Jacobi steps to reduce the error by 25% while Gauss-Seidel does the same job in one step.

This rule of thumb — “each Gauss-Seidel step is worth two Jacobi steps” — holds for many examples, in particular all convergent two-dimensional systems. However it is possible to construct larger systems for which Jacobi converges and Gauss-Seidel does not (or vice versa).

One final point worth emphasising is that it may not be immediately obvious how to order the equations. Apart from ensuring that the diagonal entries are all nonzero, there appears to be quite a bit of freedom (there are $n!$ ways of arranging n equations).

Example 7.5 The system from example 7.1 can be rearranged to

$$\begin{aligned} -x + 2y &= 0, \\ 2x - y &= 3, \end{aligned}$$

which suggests the Jacobi scheme

$$\begin{aligned} x_{k+1} &= 2y_k, \\ y_{k+1} &= 2x_k - 3. \end{aligned}$$

The iteration matrix for this scheme is

$$\mathbf{S}^{-1}\mathbf{T} = \begin{pmatrix} 0 & 2 \\ 2 & 0 \end{pmatrix},$$

whose eigenvalues are $\lambda = \pm 2$. Thus $\rho = 2 > 1$ and, with the equations ordered this way, Jacobi's method diverges. It is straightforward to show that Gauss-Seidel does too.

The situation demonstrated by example 7.5 is quite general. For two-dimensional systems it is usually the case that one of the obvious Jacobi schemes converges while the other diverges, and the corresponding Gauss-Seidel schemes behave likewise. For larger systems it is not obvious in advance whether or not a particular Jacobi or Gauss-Seidel scheme will converge, or what permutation of the equations gives the best convergence. However it is clearly a good idea to try to *make the diagonal be the largest entry in each row* — this will make $\mathbf{S}^{-1}\mathbf{T}$ as small as possible.

In practice it is usually *not* feasible to compute the eigenvalues of $\mathbf{S}^{-1}\mathbf{T}$ when the system is large. Therefore some trial and error may be required to formulate a scheme like Jacobi or Gauss-Seidel in such a way that it converges. However, for many classes of matrices which crop up particularly frequently (*e.g.* in the numerical solution of differential equations) the spectral radii of the Jacobi and Gauss-Seidel iteration matrices are known, and can be shown to be smaller than one.

7.7 Successive over-relaxation (SOR)

The origin of the SOR method is the observation that Gauss-Seidel converges monotonically: the approximations \mathbf{x}_k stay on the same side of the solution \mathbf{x} as k increases. This suggests going slightly further than Gauss-Seidel tells us to at each step. We introduce an *over-relaxation* or *acceleration parameter* ω : when $\omega = 1$ SOR is identical to Gauss-Seidel and acceleration is achieved by increasing ω . The best choice of ω depends on the problem being solved, but it always lies between 1 and 2.

For SOR the matrix splitting is

$$(\mathbf{D} + \omega\mathbf{L})\mathbf{x}_{k+1} = [(1 - \omega)\mathbf{D} - \omega\mathbf{U}]\mathbf{x}_k + \omega\mathbf{b}. \quad (7.4)$$

Notice that

1. if the method converges, *i.e.* $\mathbf{x}_k \rightarrow \mathbf{x}$ as $k \rightarrow \infty$, then \mathbf{x} satisfies the original problem $\mathbf{Ax} = (\mathbf{L} + \mathbf{D} + \mathbf{U})\mathbf{x} = \mathbf{b}$;

2. when $\omega = 1$, (7.4) is identical to the splitting for Gauss-Seidel.

The iteration matrix is

$$\mathbf{M}_\omega = (\mathbf{D} + \omega\mathbf{L})^{-1}[(1 - \omega)\mathbf{D} - \omega\mathbf{U}],$$

and the point of SOR is to discover the value of ω which minimises the spectral radius of \mathbf{M}_ω .

Example 7.6 For the same system as in example 7.1,

$$\mathbf{A} = \begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix} \Rightarrow \mathbf{L} = \begin{pmatrix} 0 & 0 \\ -1 & 0 \end{pmatrix}, \quad \mathbf{D} = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}, \quad \mathbf{U} = \begin{pmatrix} 0 & -1 \\ 0 & 0 \end{pmatrix}.$$

Thus the SOR iteration matrix is

$$\begin{aligned} \mathbf{M}_\omega &= \begin{pmatrix} 2 & 0 \\ -\omega & 2 \end{pmatrix}^{-1} \begin{pmatrix} 2(1-\omega) & \omega \\ 0 & 2(1-\omega) \end{pmatrix} = \frac{1}{4} \begin{pmatrix} 2 & 0 \\ \omega & 2 \end{pmatrix} \begin{pmatrix} 2(1-\omega) & \omega \\ 0 & 2(1-\omega) \end{pmatrix} \\ &= \begin{pmatrix} 1-\omega & \omega/2 \\ \omega(1-\omega)/2 & 1-\omega + \omega^2/4 \end{pmatrix}. \end{aligned}$$

The eigenvalues of \mathbf{M}_ω are given by

$$\begin{vmatrix} 1-\omega-\lambda & \omega/2 \\ \omega(1-\omega)/2 & 1-\omega + \omega^2/4 - \lambda \end{vmatrix} = 0$$

and thus

$$\lambda^2 - (2 - 2\omega + \omega^2/4)\lambda + (\omega - 1)^2 = 0. \quad (7.5)$$

When $\omega = 1$ this gives $\lambda = 0$ and $1/4$ as in example 7.4. As ω is increased, the smaller eigenvalue increases while the larger one decreases (their product is $(\omega - 1)^2$) and it can be shown that the optimal value of ω occurs when the two eigenvalues are equal.

Therefore set $\lambda_1 = \lambda_2 = \lambda$. From the quadratic equation (7.5) we can read off the sum and product of the eigenvalues:

$$\lambda_1\lambda_2 = \lambda^2 = (\omega - 1)^2 \Rightarrow \lambda = \omega - 1,$$

and

$$\begin{aligned} \lambda_1 + \lambda_2 = 2\lambda &= 2 - 2\omega + \omega^2/4 \Rightarrow 2(\omega - 1) = 2 - 2\omega + \omega^2/4 \\ &\Rightarrow \omega^2 - 16\omega + 16 = 0. \end{aligned}$$

The roots of this quadratic are $\omega = 4(2 \pm \sqrt{3})$, of which we take the one between 1 and 2:

$$\omega = 4(2 - \sqrt{3}) \approx 1.072$$

is the optimal value of the acceleration parameter for this system.

At this value of ω the two equal eigenvalues of \mathbf{M}_ω are

$$\lambda = \omega - 1 \approx 0.072,$$

and thus the spectral radius

$$\rho \approx 0.072.$$

Notice the improvement on Gauss-Seidel, which had $\rho = 1/4$. For this example SOR is nearly twice as good again as Gauss-Seidel (since $(1/4)^2 = 0.0625$ is only slightly less than 0.072).

The observation that the optimal value of ω occurs when two eigenvalues coincide holds for other systems also. However for large systems it is not practicable to find the eigenvalues explicitly and thus calculate the optimum value of ω as shown in example 7.6. Instead the best choice of ω must often be determined by trial and error, but it's certainly worth the effort. For large systems the improvement in convergence compared with Gauss-Seidel is typically much more dramatic than that seen in the two-dimensional example 7.6: each SOR step could be worth thirty or more Gauss-Seidel steps!

7.8 Implementation of SOR

In general, if \mathbf{A} has elements a_{ij} then the SOR method is

$$\begin{pmatrix} a_{11} & 0 & 0 & \dots \\ \omega a_{21} & a_{22} & 0 & \dots \\ \omega a_{31} & \omega a_{32} & a_{33} & \dots \\ \vdots & \vdots & \vdots & \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \end{pmatrix}_{k+1} = \begin{pmatrix} \omega b_1 \\ \omega b_2 \\ \omega b_3 \\ \vdots \end{pmatrix} + \\ + \begin{pmatrix} (1-\omega)a_{11} & -\omega a_{12} & -\omega a_{13} & \dots \\ 0 & (1-\omega)a_{22} & -\omega a_{23} & \dots \\ 0 & 0 & (1-\omega)a_{33} & \dots \\ \vdots & \vdots & \vdots & \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \end{pmatrix}_k.$$

By rearranging this as we did for Gauss-Seidel, we see that each step of SOR is of the form

$$(x_i)_{\text{new}} = x_i + \frac{\omega}{a_{ii}} \left[b_i - \sum_{j=1}^n a_{ij} x_j \right]. \quad (7.6)$$

Now the connection with Gauss-Seidel is much clearer: at each step we simply multiply the Gauss-Seidel correction by the acceleration parameter ω . Thus it is a simple matter to modify the code given in section 7.5 to produce an SOR code.

Exercises

7.1. Consider the linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$, where \mathbf{A} is the matrix

$$\mathbf{A} = \begin{pmatrix} 3 & 2 & 1 \\ 2 & 3 & 2 \\ 1 & 2 & 3 \end{pmatrix}.$$

- Derive the iteration matrices for the Jacobi, Gauss-Seidel and SOR methods.
- Hence show that the Jacobi method diverges, but the Gauss-Seidel method converges.

- 7.2. (a) Write a subroutine which generates the $n \times n$ finite difference matrix \mathbf{F}_n whose elements f_{ij} are defined by

$$f_{ij} = \begin{cases} 2 & i = j, \\ -1 & i = j \pm 1, \\ 0 & \text{otherwise,} \end{cases}$$

for any value of n . For example, when $n = 4$,

$$F_4 = \begin{pmatrix} 2 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 2 \end{pmatrix}.$$

Write a program which prints out the matrix \mathbf{F}_{10} to check that your subroutine works.

- (b) Write a program (using the subroutine from 7.2a) which solves the system of equations $\mathbf{F}_n \mathbf{x} = \mathbf{b}_n$ for \mathbf{x} using the SOR method, where \mathbf{F}_n is the finite difference matrix defined above and \mathbf{b}_n is the n -dimensional vector $\mathbf{b}_n = (1, 1, \dots, 1)^T$. The dimension n , acceleration parameter ω and desired accuracy of the solution should be set by the user at run time, and the program should output the solution \mathbf{x} and the number of iterations required for convergence. Test your program with the values $n = 8$, $\omega = 1.3$, desired accuracy = 0.0001.
- (c) Modify the program written in 7.2b to tabulate the number of iterations needed to achieve an accuracy of 0.0001 versus ω , where $\omega = 1, 1.05, \dots, 1.95, 2$. Print out the table of values generated when $n = 10$, and thus estimate the optimal value of ω when $n = 10$. How much better than Gauss-Seidel is SOR in this case?

7.3. [1998/99 Question 1]

- (a) Define the spectral radius $\rho(\mathbf{M})$ of a square matrix \mathbf{M} . An iterative scheme is given by $\mathbf{x}_{n+1} = \mathbf{b} + \mathbf{M}\mathbf{x}_n$, where \mathbf{b} and \mathbf{x}_n are vectors. Give the condition under which the iterative scheme converges, in terms of $\rho(\mathbf{M})$.
- (b) Consider the system of equations

$$\begin{aligned} 2x + 3y &= 7, \\ x - y &= 1. \end{aligned} \tag{*}$$

- (i). Show that a Jacobi scheme based on writing $x = (7 - 3y)/2$, $y = x - 1$ fails to converge.
- (ii). Reformulate the problem in such a way that a Jacobi scheme does converge.

- (iii). Formulate a convergent Gauss-Seidel method for the system (*), showing that the method is convergent. Compute two iterations of the method, starting from $x = 0, y = 0$.
- (iv). Use the method obtained in (iii) to write down an S.O.R. scheme for (*) and determine the value of the acceleration parameter that gives the fastest convergence.