

Convergence Analysis of Iterative Methods

We want to answer two questions:

1. When will the Jacobi or Gauss-Seidel Methods work? That is, under what conditions will they produce a sequence of approximations $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)}, \dots$ which converge to the true solution?
2. When they do work, how fast will the convergence to the true solution be? That is, what will the *rate of convergence* be?

In general, an iterative method that finds the solution to $A\mathbf{x} = \mathbf{b}$ takes the form

$$M\mathbf{x}^{(k+1)} = N\mathbf{x}^{(k)} + \mathbf{b},$$

so that

$$\mathbf{x}^{(k+1)} = M^{-1}N\mathbf{x}^{(k)} + M^{-1}\mathbf{b},$$

which we can rewrite as

$$\mathbf{x}^{(k+1)} = B\mathbf{x}^{(k)} + \tilde{\mathbf{b}},$$

where $B = M^{-1}N$ and $\tilde{\mathbf{b}} = M^{-1}\mathbf{b}$.

If the current approximation $\mathbf{x}^{(k)}$ is, in fact, the exact solution \mathbf{x} , then the iterative method should certainly find that the next iteration $\mathbf{x}^{(k+1)}$ is also the exact solution \mathbf{x} . That is, it should be that $\mathbf{x} = B\mathbf{x} + \tilde{\mathbf{b}}$, so that

$$\mathbf{x} = B\mathbf{x} + \tilde{\mathbf{b}} \Rightarrow M\mathbf{x} = N\mathbf{x} + \mathbf{b} \Rightarrow (M - N)\mathbf{x} = \mathbf{b}.$$

Of course since the original problem we are trying to solve is $A\mathbf{x} = \mathbf{b}$, it must be that M and N are chosen so that $A = M - N$. On the other hand, it turns out that choosing $A = M - N$ does not necessarily guarantee that the iterative method will find a sequence of vectors $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)}, \dots$ that converges to the true solution \mathbf{x} .

Whether or not a particular method will work will depend on the matrix $B = M^{-1}N$. In fact, it turns out that in general the matrix B *completely* determines the convergence (or not) of an iterative method. In particular, the initial guess generally has no effect on whether or not a particular method is convergent or on the rate of convergence, although if the initial guess is far away from the true solution, more iterations will be required to get an acceptable approximation for the true solution than if the initial guess were closer to the true solution.

To understand the convergence properties of an iterative method $\mathbf{x}^{(k+1)} = B\mathbf{x}^{(k)} + \tilde{\mathbf{b}}$, we subtract the equation $\mathbf{x}^{(k+1)} = B\mathbf{x}^{(k)} + \tilde{\mathbf{b}}$ from the equation $\mathbf{x} = B\mathbf{x} + \tilde{\mathbf{b}}$, which gives us $\mathbf{x} - \mathbf{x}^{(k+1)} = B(\mathbf{x} - \mathbf{x}^{(k)})$. That is, where the current error is $\mathbf{e}^{(k)} = \mathbf{x} - \mathbf{x}^{(k)}$, we have

$$\mathbf{e}^{(k)} = B\mathbf{e}^{(k-1)} = B(B\mathbf{e}^{(k-2)}) = B^2\mathbf{e}^{(k-2)} = \dots = B^k\mathbf{e}^{(0)}.$$

To be clear, the superscript of matrix B is the power of B , while the superscript of vector \mathbf{e} (inside parentheses) is the iteration number to which this particular error corresponds. We want $\mathbf{e}^{(k)} \rightarrow \mathbf{0}$ as $k \rightarrow \infty$. Since $\|\mathbf{e}^{(k)}\| = \|B^k\mathbf{e}_0\| \leq \|B\|^k \|\mathbf{e}_0\|$, then we will have $\|\mathbf{e}^{(k)}\| \rightarrow 0$

(which is the same as $e^{(k)} \rightarrow \theta$) if $\|B\| < 1$. Just like the norm of a vector, the norm $\|B\|$ of matrix B tells us the “size” of the matrix, or more precisely, how much bigger or smaller Bv will be compare to v . As we show below, there is a particular B matrix for each of the Jacobi and Gauss-Seidel Methods. For each method, the smaller $\|B\|$ is, the faster the method will converge, or if $\|B\| \geq 1$, neither method will normally converge.

One condition sometimes encountered in practice that will guarantee that $\|B\| < 1$ is that matrix A is *strictly diagonally dominant*. A matrix is strictly diagonally dominant if for each of its rows, the absolute value of the *diagonal* element is larger than the sum of the absolute values of the *off-diagonal* elements. That is,

$$\text{where } A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix} \quad \text{we have} \quad \begin{array}{l} |a_{11}| > |a_{12}| + |a_{13}| + \cdots + |a_{1n}| \\ |a_{22}| > |a_{21}| + |a_{23}| + \cdots + |a_{2n}| \\ \vdots \\ |a_{nn}| > |a_{n1}| + |a_{n2}| + \cdots + |a_{nn-1}| \end{array}$$

Example 3 Matrix A_1 is diagonally dominant and matrix A_2 is not:

$$A_1 = \begin{bmatrix} 4 & 1 & 2 \\ 2 & -5 & -2 \\ -1 & 2 & 7 \end{bmatrix} \quad \begin{array}{l} |4| > |1| + |2| \\ |-5| > |2| + |-2| \\ |7| > |-1| + |2| \end{array} \quad A_2 = \begin{bmatrix} 4 & 2 & 2 \\ 4 & -5 & -2 \\ -1 & 2 & 7 \end{bmatrix} \quad \begin{array}{l} |4| = |2| + |2| \\ |-5| < |4| + |-2| \\ |7| > |-1| + |2| \end{array}$$

Analysis of Jacobi and Gauss-Seidel Methods for 2 x 2 Systems

We will continue our discussion with only the 2 x 2 case, as the Java applet to be used in doing the homework exercises deals with this case. As discussed earlier, the Jacobi and Gauss-Seidel Methods are both of the form $x^{(k+1)} = Bx^{(k)} + \tilde{b}$, where their B matrices for the 2 x 2 case are

Method	B
Jacobi	$D^{-1}(-L-U) = \begin{bmatrix} a_{11} & 0 \\ 0 & a_{22} \end{bmatrix}^{-1} \begin{bmatrix} 0 & -a_{12} \\ -a_{21} & 0 \end{bmatrix} = \begin{bmatrix} 0 & -\frac{a_{12}}{a_{11}} \\ -\frac{a_{21}}{a_{22}} & 0 \end{bmatrix}$
Gauss-Seidel	$(L+D)^{-1}(-U) = \begin{bmatrix} a_{11} & 0 \\ a_{21} & a_{22} \end{bmatrix}^{-1} \begin{bmatrix} 0 & -a_{12} \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & -\frac{a_{12}}{a_{11}} \\ 0 & \frac{a_{12} a_{21}}{a_{11} a_{22}} \end{bmatrix}$

Notice for both methods that the diagonal elements of A must be non-zero, $a_{11} \neq 0$ and $a_{22} \neq 0$. It turns out that if $n \times n$ matrix B has a full set of n distinct eigenvectors (which is always the case for 2×2 systems using Jacobi or Gauss-Seidel Methods), then $\|B\| = |\lambda_{\max}|$, where λ_{\max} is the eigenvalue of matrix B with largest magnitude. Consequently, in the 2×2 case, the Jacobi and Gauss-Seidel Methods are guaranteed to converge if all of the eigenvalues of the matrix B corresponding to that method are of magnitude < 1 . This also includes the case that B has complex eigenvalues. We note that for $n \times n$ systems, things are more complicated.

We have now answered the first question posed at the beginning of this section. Because $\|e^{(k)}\| \leq \|B\|^k \|e_0\|$, the second question is also answered. For example, if $\|B\| = 0.5$, then the error $e^{(k)} = x - x^{(k)}$ will be cut approximately in half by each additional iteration. The size $\|B\|$ is directly proportional to the size of the eigenvalues of B . Consequently, a major goal in designing an iterative method is that the corresponding B matrix has eigenvalues that are as small as possible.

The eigenvalues and corresponding eigenvectors for the Jacobi and Gauss-Seidel Methods are

Method	Eigenvalues	Eigenvectors
Jacobi	$\lambda_1 = \sqrt{\frac{a_{12} a_{21}}{a_{11} a_{22}}}, \quad \lambda_2 = -\sqrt{\frac{a_{12} a_{21}}{a_{11} a_{22}}}$	$v_1 = \begin{bmatrix} 1 \\ -\sqrt{\frac{a_{11} a_{21}}{a_{12} a_{22}}} \end{bmatrix}, \quad v_2 = \begin{bmatrix} 1 \\ \sqrt{\frac{a_{11} a_{21}}{a_{12} a_{22}}} \end{bmatrix}$
Gauss-Seidel	$\lambda_1 = 0, \quad \lambda_2 = -\frac{a_{12} a_{21}}{a_{11} a_{22}}$	$v_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad v_2 = \begin{bmatrix} 1 \\ -\frac{a_{21}}{a_{22}} \end{bmatrix}$

Notice for both methods that we will have $\|B\| = |\lambda_{\max}| < 1$ if $\left| \frac{a_{12} a_{21}}{a_{11} a_{22}} \right| < 1$. Also notice that the non-zero eigenvalue for the Gauss-Seidel Method is the square of the two eigenvalues for the Jacobi Method. As a result, when dealing with 2×2 systems, $\|B_{GS}\| = \|B_{Jacobi}\|^2$, which means that the rate of convergence of the Gauss-Seidel Method is the square of the rate of convergence of the Jacobi Method. For example, if each iteration of the Jacobi Method causes the error to be halved, each iteration of the Gauss-Seidel Method will cause the error to be quartered, since $\frac{1}{4} = \left(\frac{1}{2}\right)^2$. Another way to look at this is that, for 2×2 systems, approximately twice as many iterations of the Jacobi Method iterations are needed to achieve the same level of accuracy (to the true solution x) as the Gauss-Seidel Method. We note that for general $n \times n$ systems, things are more complicated than (but similar to) the 2×2 case that we have been discussing.

SOR Method

A third iterative method, called the Successive Overrelaxation (SOR) Method, is an improvement on the Gauss-Seidel Method. The basic idea is this: in finding $\mathbf{x}^{(k+1)}$ given $\mathbf{x}^{(k)}$, we move a certain amount in a particular direction from $\mathbf{x}^{(k)}$ to $\mathbf{x}^{(k+1)}$. This direction is simply the vector $\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}$, since $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})$. Assuming that the direction from $\mathbf{x}^{(k)}$ to $\mathbf{x}^{(k+1)}$ is taking us closer, but not all the way, to the true solution \mathbf{x} , the basic idea of the SOR Method then is to move in the same direction $\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}$, but farther along that direction.

To derive the SOR Method from the Gauss-Seidel Method, notice that we could also write the Gauss-Seidel equation as

$$D\mathbf{x}^{(k+1)} = \mathbf{b} - L\mathbf{x}^{(k+1)} - U\mathbf{x}^{(k)}$$

so that

$$\mathbf{x}^{(k+1)} = D^{-1}[\mathbf{b} - L\mathbf{x}^{(k+1)} - U\mathbf{x}^{(k)}].$$

We can subtract $\mathbf{x}^{(k)}$ from both sides to get

$$\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} = D^{-1}[\mathbf{b} - L\mathbf{x}^{(k+1)} - D\mathbf{x}^{(k)} - U\mathbf{x}^{(k)}].$$

One can think of this as the *Gauss-Seidel correction* $(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})_{GS}$. As suggested above, it turns out that convergence $\mathbf{x}^{(k)} \rightarrow \mathbf{x}$ is often faster if we go beyond the Gauss-Seidel correction. The idea of the SOR Method is to iterate

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \omega(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})_{GS}$$

where

$$(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})_{GS} = D^{-1}[\mathbf{b} - L\mathbf{x}^{(k+1)} - D\mathbf{x}^{(k)} - U\mathbf{x}^{(k)}],$$

and where generally $1 < \omega < 2$. Notice that if $\omega = 1$, then we simply have the Gauss-Seidel Method. Written out in detail, the SOR Method is

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

We can multiply both sides by matrix D and divide both sides by ω , and rewrite this as

$$\frac{1}{\omega} D\mathbf{x}^{(k+1)} = \frac{1}{\omega} D\mathbf{x}^{(k)} + [\mathbf{b} - L\mathbf{x}^{(k+1)} - D\mathbf{x}^{(k)} - U\mathbf{x}^{(k)}],$$

then collect the $\mathbf{x}^{(k+1)}$ terms on the left hand side to get

$$\begin{aligned} (L + \frac{1}{\omega} D)\mathbf{x}^{(k+1)} &= \frac{1}{\omega} D\mathbf{x}^{(k)} + [\mathbf{b} - D\mathbf{x}^{(k)} - U\mathbf{x}^{(k)}] \\ \Rightarrow (L + \frac{1}{\omega} D)\mathbf{x}^{(k+1)} &= (\frac{1}{\omega} D - D - U)\mathbf{x}^{(k)} + \mathbf{b} \\ \Rightarrow \mathbf{x}^{(k+1)} &= (L + \frac{1}{\omega} D)^{-1}[(\frac{1}{\omega} D - D - U)\mathbf{x}^{(k)} + \mathbf{b}] \end{aligned}$$

Notice that the SOR Method is also of the form $\mathbf{x} = B\mathbf{x} + \tilde{\mathbf{b}}$, so that the convergence analysis just done for the Jacobi and Gauss-Seidel Methods can also be done for the SOR Method. The B matrix which determines the convergence of the SOR Method is $(L + \frac{1}{\omega} D)^{-1}(\frac{1}{\omega} D - D - U)$. From this point of view, the idea is to choose a value of ω which minimizes

$\|(L + \frac{1}{\omega}D)^{-1}(\frac{1}{\omega}D - D - U)\|$. As we did earlier for the Jacobi and Gauss-Seidel Methods, it is possible to find the eigenvalues and eigenvectors for the B matrix when using the SOR Method for 2×2 systems. However, because it is significantly more complicated, we do not do the derivations here. It is left as an exercise for the ambitious student (or the challenging instructor).

Algorithms

In practice, we are usually using a computer to do the iterations, so we need implementable algorithms in order to use the above methods for $n \times n$ systems. We conclude by giving these algorithms for finding element $x_i^{(k+1)}$, given $x_1^{(k)}, x_2^{(k)}, \dots, x_n^{(k)}$. We note that although these particular algorithms are not quite optimally efficient, writing the algorithms this way makes more obvious the slight (but important) differences between the three methods.

Method	Algorithm for performing iteration $k + 1$: For $i = 1$ to n do:
Jacobi	$x_i^{(k+1)} = x_i^{(k)} + \frac{1}{a_{ii}} \left(b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} \right)$
Gauss-Seidel	$x_i^{(k+1)} = x_i^{(k)} + \frac{1}{a_{ii}} \left(b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} \right)$
SOR	$x_i^{(k+1)} = x_i^{(k)} + \frac{\omega}{a_{ii}} \left(b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} \right)$