

For the **SOR method**, we have $\mathbf{Q} = 1/\omega(\mathbf{D} - \omega\mathbf{C}_L)$, so

$$\begin{aligned}\mathcal{L}_\omega &= \mathbf{I} - \mathbf{Q}^{-1}\mathbf{A} = (\mathbf{D} - \omega\mathbf{C}_L)^{-1}[\omega\mathbf{C}_U + (1 - \omega)\mathbf{D}] \\ \mathbf{h} &= \mathbf{Q}^{-1}\mathbf{b} = \omega(\mathbf{D} - \omega\mathbf{C}_L)^{-1}\mathbf{b}\end{aligned}$$

Another View of Overrelaxation

In some cases, the rate of convergence of the basic iterative scheme (2) can be improved by the introduction of an auxiliary vector and an *acceleration parameter* ω as follows:

$$\begin{aligned}\mathbf{Q}\mathbf{z}^{(k)} &= (\mathbf{Q} - \mathbf{A})\mathbf{x}^{(k-1)} + \mathbf{b} \\ \mathbf{x}^{(k)} &= \omega\mathbf{z}^{(k)} + (1 - \omega)\mathbf{x}^{(k-1)}\end{aligned}$$

or

$$\mathbf{x}^{(k)} = \omega\{(\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A})\mathbf{x}^{(k-1)} + \mathbf{Q}^{-1}\mathbf{b}\} + (1 - \omega)\mathbf{x}^{(k-1)}$$

The parameter ω gives a weighting in favor of the updated values. When $\omega = 1$, this procedure reduces to the basic iterative method, and when $1 < \omega < 2$, the rate of convergence may be improved, which is called **overrelaxation**. When $\mathbf{Q} = \mathbf{D}$, we have the **Jacobi overrelaxation (JOR) method**:

$$\mathbf{x}^{(k)} = \omega\{\mathbf{B}\mathbf{x}^{(k-1)} + \mathbf{h}\} + (1 - \omega)\mathbf{x}^{(k-1)}$$

Overrelaxation has particular advantages when used with the Gauss-Seidel method in a slightly different way:

$$\begin{aligned}\mathbf{D}\mathbf{z}^{(k)} &= \mathbf{C}_L\mathbf{x}^{(k)} + \mathbf{C}_U\mathbf{x}^{(k-1)} + \mathbf{b} \\ \mathbf{x}^{(k)} &= \omega\mathbf{z}^{(k)} + (1 - \omega)\mathbf{x}^{(k-1)}\end{aligned}$$

and we have the **SOR method**:

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

Conjugate Gradient Method

The conjugate gradient method is one of the most popular iterative methods for solving sparse systems of linear equations. This is particularly true for systems that arise in the numerical solutions of partial differential equations.

We begin with a brief presentation of definitions and associated notation. (Some of them are presented more fully in Chapter 16.) Assume that the real $n \times n$ matrix \mathbf{A} is **symmetric**, meaning that $\mathbf{A}^T = \mathbf{A}$. The **inner product** of two vectors $\mathbf{u} = (u_1, u_2, \dots, u_n)$ and $\mathbf{v} = (v_1, v_2, \dots, v_n)$ can be written as $\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u}^T \mathbf{v} = \sum_{i=1}^n u_i v_i$, which is the scalar sum. Note that $\langle \mathbf{u}, \mathbf{v} \rangle = \langle \mathbf{v}, \mathbf{u} \rangle$. If \mathbf{u} and \mathbf{v} are mutually orthogonal, then $\langle \mathbf{u}, \mathbf{v} \rangle = 0$. An **A-inner product** of two vectors \mathbf{u} and \mathbf{v} is defined as

$$\langle \mathbf{u}, \mathbf{v} \rangle_A = \langle \mathbf{A}\mathbf{u}, \mathbf{v} \rangle = \mathbf{u}^T \mathbf{A}^T \mathbf{v}$$

Two nonzero vectors \mathbf{u} and \mathbf{v} are **A-conjugate** if $\langle \mathbf{u}, \mathbf{v} \rangle_A = 0$. An $n \times n$ matrix \mathbf{A} is **positive definite** if

$$\langle \mathbf{x}, \mathbf{x} \rangle_A > 0$$

for all nonzero vectors $\mathbf{x} \in \mathbb{R}^n$. In general, expressions such as $\langle \mathbf{u}, \mathbf{v} \rangle$ and $\langle \mathbf{u}, \mathbf{v} \rangle_A$ reduce to 1×1 matrices and are treated as scalar values. A **quadratic form** is a scalar quadratic function of a vector of the form

$$f(\mathbf{x}) = \frac{1}{2} \langle \mathbf{x}, \mathbf{x} \rangle_A - \langle \mathbf{b}, \mathbf{x} \rangle + c$$

Here, \mathbf{A} is a matrix, \mathbf{x} and \mathbf{b} are vectors, and c is a scalar constant. The **gradient** of a quadratic form

$$f'(\mathbf{x}) = \left[\partial f(\mathbf{x}) / \partial x_1, \quad \partial f(\mathbf{x}) / \partial x_2, \quad \dots, \quad \partial f(\mathbf{x}) / \partial x_n \right]^T$$

We can derive the following:

$$f'(\mathbf{x}) = \frac{1}{2} \mathbf{A}^T \mathbf{x} + \frac{1}{2} \mathbf{A} \mathbf{x} - \mathbf{b}$$

If \mathbf{A} is symmetric, this reduces to

$$f'(\mathbf{x}) = \mathbf{A} \mathbf{x} - \mathbf{b}$$

Setting the gradient to zero, we obtain the linear system to be solved, $\mathbf{A} \mathbf{x} = \mathbf{b}$. Therefore, the solution of $\mathbf{A} \mathbf{x} = \mathbf{b}$ is a critical point of $f(\mathbf{x})$. If \mathbf{A} is symmetric and positive definite, then $f(\mathbf{x})$ is minimized by the solution of $\mathbf{A} \mathbf{x} = \mathbf{b}$. So an alternative way of solving the linear system $\mathbf{A} \mathbf{x} = \mathbf{b}$ is by finding an \mathbf{x} that minimizes $f(\mathbf{x})$.

We want to solve the linear system

$$\mathbf{A} \mathbf{x} = \mathbf{b}$$

where the $n \times n$ matrix \mathbf{A} is symmetric and positive definite.

Suppose that $\{\mathbf{p}^{(1)}, \mathbf{p}^{(2)}, \dots, \mathbf{p}^{(k)}, \dots, \mathbf{p}^{(n)}\}$ is a set containing a sequence of n mutually conjugate **direction vectors**. Then they form a basis for the space \mathbb{R}^n . Hence, we can expand the true solution vector \mathbf{x}^* of $\mathbf{A} \mathbf{x} = \mathbf{b}$ into a linear combination of these basis vectors:

$$\mathbf{x}^* = \alpha_1 \mathbf{p}^{(1)} + \alpha_2 \mathbf{p}^{(2)} + \dots + \alpha^{(k)} \mathbf{p}^{(k)} + \dots + \alpha_n \mathbf{p}^{(n)}$$

where the coefficients are given by

$$\alpha_k = \langle \mathbf{p}^{(k)}, \mathbf{b} \rangle / \langle \mathbf{p}^{(k)}, \mathbf{p}^{(k)} \rangle_A$$

This can be viewed as a direct method for solving the linear system $\mathbf{A} \mathbf{x} = \mathbf{b}$: First find the sequence of n conjugate direction vectors $\mathbf{p}^{(k)}$, and then compute the coefficients α_k . However, in practice, this approach is impractical because it would take too much computer time and storage.

On the other hand, if we view the conjugate gradient method as an iterative method, then we could solve large sparse linear systems in a reasonable amount of time and storage. The key is carefully choosing a small set of the conjugate direction vectors $\mathbf{p}^{(k)}$ so that we do not need them all to obtain a good approximation to the true solution vector.

Start with an initial guess $\mathbf{x}^{(0)}$ to the true solution \mathbf{x}^* . We can assume without loss of generality that $\mathbf{x}^{(0)}$ is the zero vector. The true solution \mathbf{x}^* is also the unique minimizer of

$$f(\mathbf{x}) = \frac{1}{2} \langle \mathbf{x}, \mathbf{x} \rangle_A - \langle \mathbf{x}, \mathbf{x} \rangle = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{x}^T \mathbf{x}$$

for $\mathbf{x} \in \mathbb{R}^n$. This suggests taking the first basis vector $\mathbf{p}^{(1)}$ to be the gradient of f at $\mathbf{x} = \mathbf{x}^{(0)}$, which equals $-\mathbf{b}$. The other vectors in the basis are now conjugate to the gradient—hence

the name *conjugate gradient method*. The k th residual vector is

$$\mathbf{r}^{(k)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}$$

The gradient descent method moves in the direction $\mathbf{r}^{(k)}$. Take the direction closest to the gradient vector $\mathbf{r}^{(k)}$ by insisting that the direction vectors $\mathbf{p}^{(k)}$ be conjugate to each other. Putting all this together, we obtain the expression

$$\mathbf{p}^{(k+1)} = \mathbf{r}^{(k)} - [\langle \mathbf{p}^{(k)}, \mathbf{r}^{(k)} \rangle_A / \langle \mathbf{p}^{(k)}, \mathbf{p}^{(k)} \rangle_A] \mathbf{p}_k$$

After some simplifications, the algorithm is obtained for solving the linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$, where the coefficient matrix \mathbf{A} is real, symmetric, and positive definite. The input vector $\mathbf{x}^{(0)}$ is an initial approximation to the solution or the zero vector.

In theory, the conjugate gradient iterative method solves a system of n linear equations in at most n steps, if the matrix \mathbf{A} is symmetric and positive definite. Moreover, the n th iterative vector $\mathbf{x}^{(n)}$ is the unique minimizer of the quadratic function $q(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{A}\mathbf{x} - \mathbf{x}^T \mathbf{b}$. When the conjugate gradient method was introduced by Hestenes and Stiefel [1952], the initial interest in it waned once it was discovered that this finite-termination property was not obtained in practice. But two decades later, there was renewed interest in this method when it was viewed as an iterative process by Reid [1971] and others. In practice, the solution of a system of linear equations can often be found with satisfactory precision in a number of steps considerably less than the order of the system.

Here is a pseudocode for the **conjugate gradient algorithm**:

```

 $k \leftarrow 0$ ;  $\mathbf{x} \leftarrow \mathbf{0}$ ;  $\mathbf{r} \leftarrow \mathbf{b} - \mathbf{A}\mathbf{x}$ ;  $\delta \leftarrow \langle \mathbf{r}, \mathbf{r} \rangle$ 
while ( $\sqrt{\delta} > \varepsilon \sqrt{\langle \mathbf{b}, \mathbf{b} \rangle}$ ) and ( $k < k_{\max}$ )
     $k \leftarrow k + 1$ 
    if  $k = 1$  then
         $\mathbf{p} \leftarrow \mathbf{r}$ 
    else
         $\beta \leftarrow \delta / \delta_{\text{old}}$ 
         $\mathbf{p} \leftarrow \mathbf{r} + \beta \mathbf{p}$ 
    end if
     $\mathbf{w} \leftarrow \mathbf{A}\mathbf{p}$ 
     $\alpha \leftarrow \delta / \langle \mathbf{p}, \mathbf{w} \rangle$ 
     $\mathbf{x} \leftarrow \mathbf{x} + \alpha \mathbf{p}$ 
     $\mathbf{r} \leftarrow \mathbf{r} - \alpha \mathbf{w}$ 
     $\delta_{\text{old}} \leftarrow \delta$ 
     $\delta \leftarrow \langle \mathbf{r}, \mathbf{r} \rangle$ 
end while

```

Here, ε is a parameter used in the convergence criterion (such as $\varepsilon = 10^{-5}$), and k_{\max} is the maximum number of iterations allowed. Usually, the number of iterations needed is much less than the size of the linear system. We save the previous value of δ in the variable δ_{old} . If a good guess for the solution vector \mathbf{x} is known, then it should be used as an initial vector instead of zero. The variable ε is the desired convergence tolerance. The algorithm produces not only a sequence of vectors $\mathbf{x}^{(i)}$ that converges to the solution but an orthogonal sequence of residual vectors $\mathbf{r}^{(i)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(i)}$ and an \mathbf{A} -orthogonal sequence of

search direction vectors $\mathbf{p}^{(i)}$, namely, $\langle \mathbf{r}^{(i)}, \mathbf{r}^{(j)} \rangle = 0$ if $i \neq j$ and $\langle \mathbf{p}^{(i)}, \mathbf{A}\mathbf{p}^{(j)} \rangle = 0$ if $i \neq j$. (The main computational features of the conjugate gradient algorithm are complicated to derive, but the final conclusion is that in each step, only *one* matrix-vector multiplication is required and only a few dot-products are computed. These are extremely desirable attributes in solving large and sparse linear systems. Also, unlike Gaussian elimination, there is no fill-in, so only the nonzero entries in \mathbf{A} need to be stored in the computer memory. For some partial differential equation problems, the equations in the linear system can be represented by stencils that describe the nonzero structure within the coefficient matrix. Sometimes these stencils are used in a computer program rather than storing the nonzero entries in the coefficient matrix.

EXAMPLE 5 Use the conjugate gradient method to solve this linear system:

$$\begin{bmatrix} 2 & -1 & 0 \\ -1 & 3 & -1 \\ 0 & -1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 8 \\ -5 \end{bmatrix}$$

Solution Programming the pseudocode, we obtain the iterates

$$\begin{aligned} \mathbf{x}^{(0)} &= [0.00000, 0.00000, 0.00000]^T \\ \mathbf{x}^{(1)} &= [0.29221, 2.33766, -1.46108]^T \\ \mathbf{x}^{(2)} &= [1.82254, 2.60772, -1.55106]^T \\ \mathbf{x}^{(3)} &= [2.00000, 3.00000, -1.00000]^T \end{aligned}$$

In only three iterations, we have the answer accurate to full machine precision, which illustrates the finite termination property. The matrix \mathbf{A} is symmetric positive definite and the eigenvalues of \mathbf{A} are 1, 2, 4. This simple example may be a bit misleading because one cannot expect such rapid convergence in realistic applications. (The rate of convergence depends on various properties of the linear system.) In fact, the above example is too small to illustrate the power of advanced iterative methods on very large and sparse systems. ■

The conjugate gradient method may converge slowly when the matrix \mathbf{A} is ill-conditioned; however, the convergence can be accelerated by a technique called **preconditioning**. This involves a matrix \mathbf{M}^{-1} that approximates \mathbf{A} so that $\mathbf{M}^{-1}\mathbf{A}$ is well-conditioned and $\mathbf{M}\mathbf{x} = \mathbf{y}$ is easily solved. For many very large and sparse linear systems, preconditioned conjugate gradient methods have now become the iterative methods of choice! For additional details, see Golub and Van Loan [1996] as well as many other standard textbooks and references.

Summary

(1) For the linear system

$$\mathbf{A}\mathbf{x} = \mathbf{b}$$

the general form of an iterative method is

$$\mathbf{x}^{(k)} = \mathcal{G}\mathbf{x}^{(k-1)} + \mathbf{h}$$