

Mathematical Tripos Part II: Michaelmas Term 2024

Numerical Analysis – Lecture 18

Approach 4.20 (Minimization of quadratic function) The methods we considered so far for solving $Ax = b$, namely Jacobi, Gauss-Seidel, and those with relaxation, fit into the scheme

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + c_k \mathbf{d}^{(k)},$$

where we were aimed at getting $\rho(H) < 1$ for the iteration matrix H . Say, for Jacobi with relaxation, we set $c_k = \omega$ and $\mathbf{d}^{(k)} = D^{-1}(\mathbf{b} - A\mathbf{x}^{(k)})$.

For solving $Ax = b$ with a (positive definite) matrix $A > 0$, there is a different approach to constructing good iterative methods. It is based on successive minimization of the quadratic function

$$F(\mathbf{x}^{(k)}) := \|\mathbf{x}^* - \mathbf{x}^{(k)}\|_A^2 = \|\mathbf{e}^{(k)}\|_A^2,$$

since the minimizer is clearly the exact solution. Here, $\|\mathbf{y}\|_A := (A\mathbf{y}, \mathbf{y})^{1/2} := \sqrt{\mathbf{y}^T A \mathbf{y}}$ is a Euclidean-type distance which is well-defined for $A > 0$. So, at each step k , we are decreasing the A -distance between $\mathbf{x}^{(k)}$ and the exact solution \mathbf{x}^* . Thus, for a symmetric positive definite $A > 0$, we choose an iterative method that provides the descent condition

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + c_k \mathbf{d}^{(k)} \Rightarrow F(\mathbf{x}^{(k+1)}) < F(\mathbf{x}^{(k)}). \quad (4.5)$$

An equivalent approach is to minimize the quadratic function

$$F_1(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T A \mathbf{x} - \mathbf{x}^T \mathbf{b},$$

which attains its minimum when $\nabla F_1(\mathbf{x}) = A\mathbf{x} - \mathbf{b} = 0$, and which does not involve the unknown \mathbf{x}^* . It is easy to check that $F_1(\mathbf{x}) = \frac{1}{2} F(\mathbf{x}) - \frac{1}{2} c$, where $c = \mathbf{x}^{*T} A \mathbf{x}^*$ is a constant independent of k , hence equivalence.

Example 4.21 Both the Jacobi and the Gauss-Seidel methods satisfy (4.5), precisely

$$(A\mathbf{e}^{(k+1)}, \mathbf{e}^{(k+1)}) = (A\mathbf{e}^{(k)}, \mathbf{e}^{(k)}) - (C\mathbf{y}^{(k)}, \mathbf{y}^{(k)}) < (A\mathbf{e}^{(k)}, \mathbf{e}^{(k)}),$$

$$\text{where for Gauss-Seidel: } C = D > 0, \quad \mathbf{y}^{(k)} := (L_0 + D)^{-1} A\mathbf{e}^{(k)};$$

$$\text{and for Jacobi: } C = 2D - A > 0, \quad \mathbf{y}^{(k)} := D^{-1} A\mathbf{e}^{(k)}.$$

Method 4.22 (A-orthogonal projection) Next, we strengthen the descent condition (4.5), namely given $\mathbf{x}^{(k)}$ and some $\mathbf{d}^{(k)}$ (called a *search direction*), we will seek $\mathbf{x}^{(k+1)}$ from the set of vectors on the line $\ell = \{\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}\}_{\alpha \in \mathbb{R}}$ such that it makes the value of $F(\mathbf{x}^{(k+1)})$ not just smaller than $F(\mathbf{x}^{(k)})$, but as small as possible (with respect to this set), namely

$$\mathbf{x}^{(k+1)} := \arg \min_{\alpha} F(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}). \quad (4.6)$$

Lemma 4.23 The minimizer in (4.6) is given by the formula

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}, \quad \alpha_k = \frac{(\mathbf{r}^{(k)}, \mathbf{d}^{(k)})}{(A\mathbf{d}^{(k)}, \mathbf{d}^{(k)})}. \quad (4.7)$$

(This choice of α_k is referred to as exact line search.)

Proof. From the definition of F , it follows that in (4.6) we should choose the point $\mathbf{x}^{(k+1)} \in \ell$ that minimizes the A -distance between \mathbf{x}^* and the points $\mathbf{y} \in \ell$. Geometrically, it is clear that the minimum occurs when $\mathbf{x}^{(k+1)}$ is the A -orthogonal projection of \mathbf{x}^* onto the line $\ell = \{\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}\}$, i.e., when

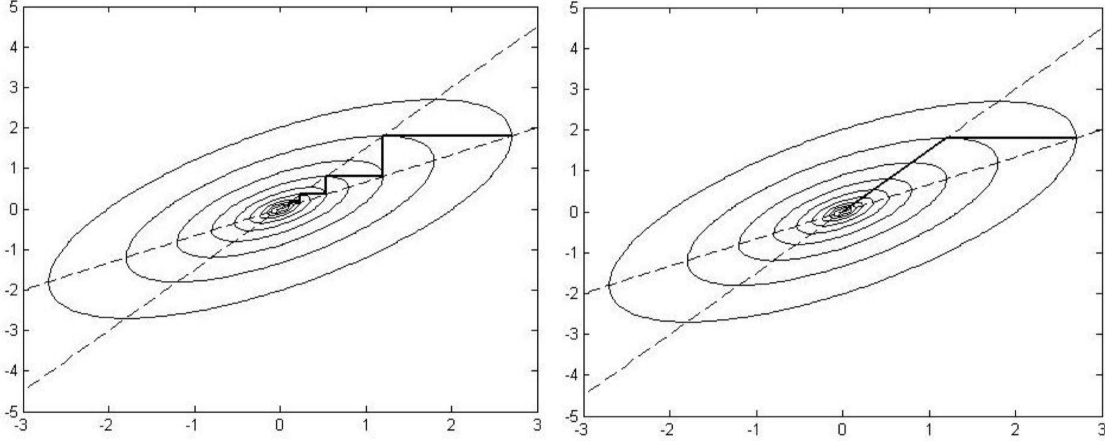
$$\mathbf{x}^* - \mathbf{x}^{(k+1)} \perp_A \mathbf{d}^{(k)} \Rightarrow A(\mathbf{x}^* - \mathbf{x}^{(k+1)}) \perp \mathbf{d}^{(k)} \Rightarrow \mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} - \alpha_k A\mathbf{d}^{(k)} \perp \mathbf{d}^{(k)}.$$

This gives expression for α_k in (4.7). □

Method 4.24 (The steepest descent method) This method takes $\mathbf{d}^{(k)} = -\nabla F_1(\mathbf{x}^{(k)}) = \mathbf{b} - A\mathbf{x}^{(k)}$ for every k , the reason being that, locally, the negative gradient of a quadratic function shows the direction of the (locally) steepest descent at a given point. Thus, the iterations have the form

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k(\mathbf{b} - A\mathbf{x}^{(k)}), \quad k \geq 0. \quad (4.8)$$

It can be proved that the sequence $(\mathbf{x}^{(k)})$ converges to the solution \mathbf{x}^* of the system $A\mathbf{x} = \mathbf{b}$ as required, but usually the speed of convergence is rather slow. The reason is that the iteration (4.8) decreases the value of $F(\mathbf{x}^{(k+1)})$ locally, relatively to $F(\mathbf{x}^{(k)})$, but the global decrease, with respect to $F(\mathbf{x}^{(0)})$, is often not that large. The use of *conjugate directions* provides a method with a global minimization property.



(a) Worst case scenario of steepest descent

(b) Conjugate gradient method applied to the same problem as in (a)

Conjugate directions Let's revisit equation (4.7) for a general direction \mathbf{d} (i.e., not necessarily equal to the negative gradient). Assume $\mathbf{x} = \mathbf{x}^{(k)}$, and let $\mathbf{e}^{(k)} = \mathbf{x}^* - \mathbf{x}^{(k)}$ be the error and $\mathbf{r}^{(k)} = \mathbf{b} - A\mathbf{x}^{(k)} = A\mathbf{e}^{(k)}$ be the residual. Then we can write $\langle \mathbf{r}^{(k)}, \mathbf{d} \rangle = \langle \mathbf{e}^{(k)}, \mathbf{d} \rangle_A$, and so for a general search direction \mathbf{d} with an exact line search, the iterate takes the form $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \frac{\langle \mathbf{e}^{(k)}, \mathbf{d} \rangle_A}{\langle \mathbf{d}, \mathbf{d} \rangle_A} \mathbf{d}$. By subtracting \mathbf{x}^* , the iterates in terms of the error $\mathbf{e}^{(k+1)}$ are given by:

$$\mathbf{e}^{(k+1)} = \mathbf{e}^{(k)} - \frac{\langle \mathbf{e}^{(k)}, \mathbf{d} \rangle_A}{\langle \mathbf{d}, \mathbf{d} \rangle_A} \mathbf{d}. \quad (4.9)$$

Geometrically, this means that $\mathbf{e}^{(k+1)}$ is the projection of $\mathbf{e}^{(k)}$ onto the hyperplane that is A -orthogonal to \mathbf{d} , i.e., we have

$$\langle \mathbf{e}^{(k+1)}, \mathbf{d} \rangle_A = 0. \quad (4.10)$$

Definition 4.25 (Conjugate directions) The vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ are *conjugate* with respect to a symmetric positive definite matrix A if they are nonzero and A -orthogonal: $\langle \mathbf{u}, \mathbf{v} \rangle_A := \langle \mathbf{u}, A\mathbf{v} \rangle = 0$.

The observation above allows us to prove the following important result.

Theorem 4.26 Let $\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots, \mathbf{d}^{(n-1)}$ be n nonzero pairwise conjugate directions, and consider the sequence of iterates

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}, \quad \alpha_k = \frac{\langle \mathbf{r}^{(k)}, \mathbf{d}^{(k)} \rangle}{\langle \mathbf{d}^{(k)}, A\mathbf{d}^{(k)} \rangle}.$$

Let $\mathbf{r}^{(k)} = \mathbf{b} - A\mathbf{x}^{(k)}$ be the residual. Then for each $k = 1, \dots, n$, $\mathbf{r}^{(k)}$ is orthogonal to $\text{span}\{\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k-1)}\}$. In particular $\mathbf{r}^{(n)} = 0$.

Proof. Since $\mathbf{r}^{(k)} = A\mathbf{e}^{(k)}$, it suffices to show that $\mathbf{e}^{(k)}$ is A -orthogonal to $\text{span}\{\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k-1)}\}$. The proof is by induction on k . For $k = 0$ there is nothing to prove. Assume the statement is true for $k \geq 0$, and consider the equation (4.9) (with $\mathbf{d} = \mathbf{d}^{(k)}$). From the induction hypothesis, and the fact that the $\mathbf{d}^{(i)}$ are pairwise conjugate directions, we see that $\mathbf{e}^{(k+1)}$ is A -orthogonal to $\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k-1)}$. Furthermore, we have already seen in (4.10) that $\langle \mathbf{e}^{(k+1)}, \mathbf{d}^{(k)} \rangle_A = 0$. Thus this shows that $\mathbf{e}^{(k+1)}$ is A -orthogonal to $\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k)}$ as desired. \square

So, if a sequence $(\mathbf{d}^{(k)})$ of conjugate directions is at hands, we have an iterative procedure with good approximation properties.

The (A -orthogonal) basis of conjugate directions is constructed by A -orthogonalization of the sequence $\{\mathbf{r}_0, A\mathbf{r}_0, A^2\mathbf{r}_0, \dots, A^{n-1}\mathbf{r}_0\}$ with $\mathbf{r}_0 = \mathbf{b} - A\mathbf{x}_0$. This is done in the way similar to orthogonalization of the monomial sequence $\{1, x, x^2, \dots, x^{n-1}\}$ using a recurrence relation.

Remark 4.27 It is possible to extend the methods for solving $A\mathbf{x} = \mathbf{b}$ with symmetric positive definite A to any other matrices by a simple trick. Suppose we want to solve $B\mathbf{x} = \mathbf{c}$, where $B \in \mathbb{R}^{n \times n}$ is nonsingular. We can convert the above system to the symmetric and positive definite setting by defining $A = B^T B$, $\mathbf{b} = B^T \mathbf{c}$ and then solving $A\mathbf{x} = \mathbf{b}$ with the conjugate gradient algorithm (or any other method for positive definite A).