

## Mathematical Tripos Part II: Michaelmas Term 2020

### Numerical Analysis – Lecture 19

**Approach 4.20 (Minimization of quadratic function)** The methods we considered so far for solving  $A\mathbf{x} = \mathbf{b}$ , 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  $A\mathbf{x} = \mathbf{b}$  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

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

**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

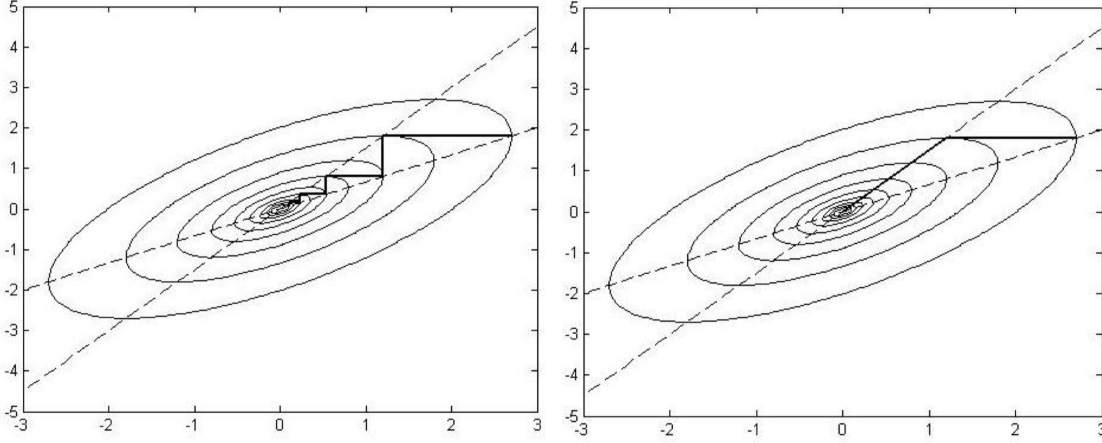
$$\mathbf{x}^* - \mathbf{x}^{(k+1)} \perp_A \mathbf{d}^{(k)} \Rightarrow A(\mathbf{x}^* - \mathbf{x}^{(k+1)}) \perp \mathbf{d}^{(k)} \Rightarrow \mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} - \alpha_k A\mathbf{d}^{(k)} \perp \mathbf{d}^{(k)}.$$

This gives expression for  $\alpha_k$  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)

**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:  $(\mathbf{u}, \mathbf{v})_A := (A\mathbf{u}, \mathbf{v}) = 0$ .

**Theorem 4.26 (Non-examinable)** Given  $A \in \mathbb{R}^{n \times n}$ ,  $A > 0$ , let  $\{\mathbf{d}^{(k)}\}_{k=0}^{n-1}$  be a set of the conjugate directions, i.e.,  $(A\mathbf{d}^{(k)}, \mathbf{d}^{(i)}) = 0$  for  $i < k$ . Then the value of  $F(\mathbf{x}^{(m+1)})$  obtained through step-by-step minimization for each  $k = 0..m$  as described in (4.7) coincides with the minimum of  $F(\mathbf{y})$  taken over all  $\mathbf{y} = \mathbf{x}^{(0)} + \sum_{k=0}^m c_k \mathbf{d}^{(k)}$  simultaneously, namely

$$\arg \min_{c_0, \dots, c_m} F(\mathbf{y}) = \mathbf{x}^{(m+1)} = \mathbf{x}^{(0)} + \sum_{k=0}^m \alpha_k \mathbf{d}^{(k)}.$$

**Proof.** Again, it is clear geometrically that the minimal  $A$ -distance between the exact solution  $\mathbf{x}^*$  and the points  $\mathbf{y}$  on the plane  $\mathcal{P} := \{\mathbf{x}^{(0)} + \sum_{k=0}^m c_k \mathbf{d}^{(k)} : c_k \in \mathbb{R}\}$  is attained when  $\mathbf{x}^{(m+1)} \in \mathcal{P}$  is the  $A$ -orthogonal projection of  $\mathbf{x}^*$  onto  $\mathcal{P}$ , i.e.,

$$\arg \min_{\mathbf{y} \in \mathcal{P}} F(\mathbf{y}) = \mathbf{x}^{(m+1)} \Leftrightarrow \mathbf{x}^* - \mathbf{x}^{(m+1)} \perp_A \{\mathbf{d}^{(k)}\}_{k=0}^m.$$

It can be shown then, that (for conjugate  $\{\mathbf{d}^{(k)}\}$ ) the latter conditions provide expressions for  $\alpha_k$  as given in (4.7).  $\square$

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  $Ax = b$  with symmetric positive definite  $A$  to any other matrices by a simple trick. Suppose we want to solve  $Bx = c$ , 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  $A = B^T B$ ,  $b = B^T c$  and then solving  $Ax = b$  with the conjugate gradient algorithm (or any other method for positive definite  $A$ ).