

# Numerical Analysis - Part II

Anders C. Hansen

Lecture 18

---

*Iterative methods for linear algebraic systems*

# Solving linear systems with iterative methods

The general *iterative* method for solving  $A\mathbf{x} = \mathbf{b}$  is a rule  $\mathbf{x}^{k+1} = f_k(\mathbf{x}^0, \mathbf{x}^1, \dots, \mathbf{x}^k)$ . We will consider the simplest ones: *linear, one-step, stationary* iterative schemes:

$$\mathbf{x}^{k+1} = H\mathbf{x}^k + \mathbf{v}, \quad \mathbf{x}^0, \mathbf{v} \in \mathbb{R}^n. \quad (1)$$

Here one chooses  $H$  and  $\mathbf{v}$  so that  $\mathbf{x}^*$ , a solution of  $A\mathbf{x} = \mathbf{b}$ , satisfies  $\mathbf{x}^* = H\mathbf{x}^* + \mathbf{v}$ , i.e. it is the fixed point of the iteration (1) (if the scheme converges). Standard terminology:

- ▶ the *iteration matrix*  $H$ ,
- ▶ the *error*  $\mathbf{e}^k := \mathbf{x}^* - \mathbf{x}^k$ ,
- ▶ the *residual*  $\mathbf{r}^k := A\mathbf{e}^k = \mathbf{b} - A\mathbf{x}^k$ .

# Linear systems in elliptic PDEs

---

As we have seen in the previous sections linear systems  $A\mathbf{x} = \mathbf{b}$ , where  $A$  is a real symmetric positive (negative) definite matrix, frequently occur in numerical methods for solving elliptic partial differential equations.

# Poisson's equation on a square

A typical example we already encountered is Poisson's equation on a square where the *five-point formula* approximation yields an  $n \times n$  system of linear equations with  $n = m^2$  unknowns  $u_{p,q}$ :

$$u_{p-1,q} + u_{p+1,q} + u_{p,q-1} + u_{p,q+1} - 4u_{p,q} = h^2 f(ph, qh) \quad (2)$$

In the *natural ordering*, when the grid points are arranged by columns,  $A$  is the following block tridiagonal matrix:

$$A = \begin{bmatrix} B & I & & & \\ I & B & I & & \\ & \ddots & \ddots & \ddots & \\ & & & I & B & I \\ & & & & I & B \end{bmatrix}, \quad B = \begin{bmatrix} -4 & 1 & & & \\ & 1 & -4 & 1 & \\ & & \ddots & \ddots & \ddots \\ & & & 1 & -4 & 1 \\ & & & & 1 & -4 \end{bmatrix}. \quad (3)$$

# The matrix $A$ is symmetric and negative definite

---

## Lemma 1

*For any ordering of the grid points, the matrix  $A$  of the system (2) is symmetric and negative definite.*

# Poisson's equation on a square

---

Note that when  $p$  or  $q$  is equal to 1 or  $m$ , then the values  $u_{0,q}$ ,  $u_{p,0}$  or  $u_{p,m+1}$ ,  $u_{m+1,q}$  are known boundary values and they should be moved to the right-hand side, thus leaving fewer unknowns on the left.

For any ordering of the grid points  $(ph, qh)$  we have shown in Lemma 1 that the matrix  $A$  of this linear system is symmetric and negative definite.

## Corollary 2

*For the linear system (2), for any ordering of the grid, both Jacobi and Gauss-Seidel methods converge.*

**Proof.** By Lemma 1,  $A$  is symmetric and negative definite, hence convergence of Gauss-Seidel. To prove convergence of the Jacobi method, we need negative definiteness of the matrix  $2D - A$ , and that follows by the same arguments as in Lemma 1: recall that the proof operates with the modulus of the off-diagonal elements and does not depend on their sign. □



# Relaxation

---

It is often possible to improve the efficiency of the splitting method by *relaxation*. Specifically, instead of letting  $(A - B)\mathbf{x}^{(k+1)} = -B\mathbf{x}^{(k)} + \mathbf{b}$ , we let

$$(A - B)\hat{\mathbf{x}}^{(k+1)} = -B\mathbf{x}^{(k)} + \mathbf{b}, \quad \text{and then} \quad \mathbf{x}^{(k+1)} = \omega\hat{\mathbf{x}}^{(k+1)} + (1 - \omega)\mathbf{x}^{(k)}$$

where  $\omega$  is a real constant called the *relaxation parameter*.

Note that  $\omega = 1$  corresponds to the standard “unrelaxed” iteration. Good choice of  $\omega$  leads to a smaller spectral radius of the iteration matrix (compared with the “unrelaxed” method), and the smaller the spectral radius, the faster the iteration converges.

# Relaxation

To this end, let us express the relaxation iteration matrix  $H_\omega$  in terms of  $H = -(A - B)^{-1}B$ . We have

$$\begin{aligned}\widehat{\mathbf{x}}^{(k+1)} = H\mathbf{x}^{(k)} + \mathbf{v} \quad \Rightarrow \quad \mathbf{x}^{(k+1)} &= \omega\widehat{\mathbf{x}}^{(k+1)} + (1 - \omega)\mathbf{x}^{(k)} \\ &= \omega H\mathbf{x}^{(k)} + (1 - \omega)\mathbf{x}^{(k)} + \omega\mathbf{v},\end{aligned}$$

hence

$$H_\omega = \omega H + (1 - \omega)I.$$

It follows that the spectra of  $H_\omega$  and  $H$  are related by the rule  $\lambda_\omega = \omega\lambda + (1 - \omega)$ , therefore one may try to choose  $\omega \in \mathbb{R}$  to minimize

$$\rho(H_\omega) = \max \{|\omega\lambda + (1 - \omega)| : \lambda \in \sigma(H)\}.$$

In general,  $\sigma(H)$  is unknown, but often we have some information about it which can be utilized to find a "good" (rather than "best") value of  $\omega$ . For example, suppose that it is known that  $\sigma(H)$  is real and resides in the interval  $[\alpha, \beta]$  where  $-1 < \alpha < \beta < 1$ . In that case we seek  $\omega$  to minimize

$$\max \{|\omega\lambda + (1 - \omega)| : \lambda \in [\alpha, \beta]\}.$$

It is readily seen that, for a fixed  $\lambda < 1$ , the function  $f(\omega) = \omega\lambda + (1 - \omega)$  is decreasing, therefore, as  $\omega$  increases (decreases) from 1 the spectrum of  $H_\omega$  moves to the left (to the right) of the spectrum of  $H$ . It is clear that the optimal location of the spectrum  $\sigma(H_\omega)$  (or of the interval  $[\alpha_\omega, \beta_\omega]$  that contains  $\sigma(H_\omega)$ ) is the one which is centralized around the origin:

$$-\left[\omega\alpha + (1-\omega)\right] = \omega\beta + (1-\omega) \quad \Rightarrow \quad \begin{aligned} \omega_{\text{opt}} &= \frac{2}{2 - (\alpha + \beta)} \\ -\alpha_{\omega_{\text{opt}}} &= \beta_{\omega_{\text{opt}}} = \frac{\beta - \alpha}{2 - (\alpha + \beta)}. \end{aligned}$$

## 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  $A > 0$ .

# Minimization of quadratic function

---

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

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + c_k \mathbf{d}^{(k)} \quad \Rightarrow \quad F(\mathbf{x}^{(k+1)}) < F(\mathbf{x}^{(k)}). \quad (4)$$

# Minimization of quadratic function

---

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.

## Quadratic function – Jacobi and Gauss–Seidel

Both the Jacobi and the Gauss–Seidel methods satisfy (4), precisely

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

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

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



**A-orthogonal projection method:** Next, we strengthen the descent condition (4), 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

$$\mathbf{x}^{(k+1)} := \arg \min_{\alpha} F(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}). \quad (5)$$

## Lemma 3

*The minimizer in (5) is given by the formula*

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}, \quad \alpha_k = \frac{(\mathbf{r}^{(k)}, \mathbf{d}^{(k)})}{(A\mathbf{d}^{(k)}, \mathbf{d}^{(k)})}. \quad (6)$$

This choice of  $\alpha_k$  is referred to as exact line search.

# A-orthogonal projection

**Proof.** From the definition of  $F$ , it follows that in (5) 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

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

This gives expression for  $\alpha_k$  in (6). □

# The steepest descent method

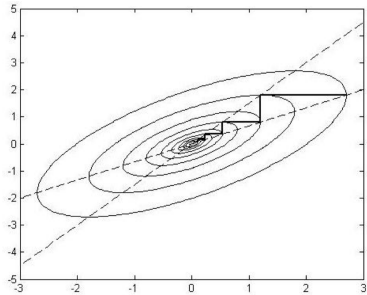
**The steepest descent method:** This method takes  $\mathbf{d}^{(k)} = -\nabla F_1(\mathbf{x}^{(k)}) = \mathbf{b} - A\mathbf{x}^{(k)}$  for every  $k$ , the reason being that, locally, the negative gradient of a quadratic function shows the direction of the (locally) steepest descent at a given point.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

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k(\mathbf{b} - A\mathbf{x}^{(k)}), \quad k \geq 0. \quad (7)$$

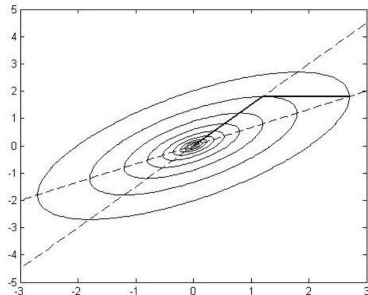
It can be proved that the sequence  $(\mathbf{x}^{(k)})$  converges to the solution  $\mathbf{x}^*$  of the system  $A\mathbf{x} = \mathbf{b}$  as required, but usually the speed of convergence is rather slow.

The reason is that the iteration (7) 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.

# Steepest descent and conjugate gradient



(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 (6) 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:

$$\mathbf{e}^{(k+1)} = \mathbf{e}^{(k)} - \frac{\langle \mathbf{e}^{(k)}, \mathbf{d} \rangle_A}{\langle \mathbf{d}, \mathbf{d} \rangle_A} \mathbf{d}. \quad (8)$$

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

$$\langle \mathbf{e}^{(k+1)}, \mathbf{d} \rangle_A = 0. \quad (9)$$

## Definition 4 (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:  
 $\langle \mathbf{u}, \mathbf{v} \rangle_A := \langle \mathbf{u}, A\mathbf{v} \rangle = 0$ .

# Conjugate gradient - Warm up

## Theorem 5

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  $\mathbf{r}^{(n)} = \mathbf{0}$ .



## Conjugate gradient - Warm up

**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 (8) (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 (9) 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.  $\square$

## Conjugate gradient - Warm up

---

So, if a sequence ( $\mathbf{d}^{(k)}$ ) of conjugate directions is at hand, 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 6

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  $A$ ).

# The conjugate gradient method

Here it is.

(A) For any initial vector  $\mathbf{x}^{(0)}$ , set  $\mathbf{d}^{(0)} = \mathbf{r}^{(0)} = \mathbf{b} - A\mathbf{x}^{(0)}$ ;

(B) For  $k \geq 0$ , calculate  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}$  and the residual

$$\mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} - \alpha_k A\mathbf{d}^{(k)}, \quad \text{with}$$
$$\alpha_k := \{ \mathbf{r}^{(k+1)} \perp \mathbf{d}^{(k)} \} = \frac{(\mathbf{r}^{(k)}, \mathbf{d}^{(k)})}{(A\mathbf{d}^{(k)}, \mathbf{d}^{(k)})}, \quad k \geq 0. \quad (10)$$

(C) For the same  $k$ , the next conjugate direction is the vector

$$\mathbf{d}^{(k+1)} = \mathbf{r}^{(k+1)} + \beta_k \mathbf{d}^{(k)}, \quad \text{with}$$
$$\beta_k := \{ \mathbf{d}^{(k+1)} \perp A\mathbf{d}^{(k)} \} = -\frac{(\mathbf{r}^{(k+1)}, A\mathbf{d}^{(k)})}{(\mathbf{d}^{(k)}, A\mathbf{d}^{(k)})}, \quad k \geq 0. \quad (11)$$