



## Classical Numerical Analysis, Chapter 07

Abner J. Salgado and Steven M. Wise

asalgad1@utk.edu swise1@utk.edu  
University of Tennessee



# Chapter 07, Part 1 of 2

## Variational and Krylov Subspace Methods

## Solving Via Optimization



In this chapter, we introduce gradient-type methods for solving the system of equations  $\mathbf{Ax} = \mathbf{f}$ , where  $A \in \mathbb{C}^{n \times n}$  is Hermitian positive definite (HPD). These are iterative methods that include the steepest descent and conjugate gradient methods. We will take advantage of the fact that solving this equation is equivalent to minimizing the quadratic function

$$E_A(\mathbf{z}) = \frac{1}{2} \mathbf{z}^H \mathbf{A} \mathbf{z} - \Re(\mathbf{z}^H \mathbf{f})$$

over  $\mathbb{C}^n$ , and will utilize some simple ideas from the theory of convex optimization.



# Basic Facts about HPD Matrices



## Theorem (properties of HPD matrices)

Suppose that  $A \in \mathbb{C}^{n \times n}$  is Hermitian positive definite (HPD). Then the following are true.

- 1 The expression

$$(\mathbf{x}, \mathbf{y})_A = (A\mathbf{x}, \mathbf{y})_2 = \mathbf{y}^H A \mathbf{x}, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{C}^n,$$

defines an inner product on  $\mathbb{C}^n$ .

- 2 The object  $\|\mathbf{x}\|_A = \sqrt{\mathbf{x}^H A \mathbf{x}}$ , where  $\mathbf{x} \in \mathbb{C}^n$  defines a norm on  $\mathbb{C}^n$ .
- 3 Let the eigenvalues of  $A$  be ordered so that  $0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ . Then

$$\sqrt{\lambda_1} \|\mathbf{x}\|_2 \leq \|\mathbf{x}\|_A \leq \sqrt{\lambda_n} \|\mathbf{x}\|_2,$$

for any  $\mathbf{x} \in \mathbb{C}^n$ .

- 4 Let  $\mathbf{f} \in \mathbb{C}^n$  be given. Then  $\mathbf{x} = A^{-1}\mathbf{f}$  if and only if  $\mathbf{x}$  minimizes the quadratic function  $E_A : \mathbb{C}^n \rightarrow \mathbb{R}$  defined by

$$E_A(\mathbf{z}) = \frac{1}{2} \mathbf{z}^H A \mathbf{z} - \Re(\mathbf{z}^H \mathbf{f}).$$



## Proof of (4).

(4:  $\implies$ ): Suppose that  $\mathbf{x} = \mathbf{A}^{-1}\mathbf{f}$ . Let  $\mathbf{y} \in \mathbb{C}^n$  be arbitrary and consider

$$\begin{aligned} E_A(\mathbf{x} + \mathbf{y}) &= \frac{1}{2} (\mathbf{x} + \mathbf{y})^H \mathbf{A} (\mathbf{x} + \mathbf{y}) - \Re(\mathbf{x}^H \mathbf{f}) - \Re(\mathbf{y}^H \mathbf{f}) \\ &= \frac{1}{2} \mathbf{x}^H \mathbf{A} \mathbf{x} + \frac{1}{2} \mathbf{y}^H \mathbf{A} \mathbf{y} + \Re(\mathbf{y}^H (\mathbf{A} \mathbf{x} - \mathbf{f})) - \Re(\mathbf{x}^H \mathbf{f}) \\ &= E_A(\mathbf{x}) + \frac{1}{2} \mathbf{y}^H \mathbf{A} \mathbf{y} \\ &\geq E_A(\mathbf{x}), \end{aligned}$$

where, in the last step, we used that  $\mathbf{A}$  is HPD.

An alternate approach is to establish the following:

$$E_A(\mathbf{z}) = \frac{1}{2} \|\mathbf{z} - \mathbf{A}^{-1}\mathbf{f}\|_{\mathbf{A}}^2 - \frac{1}{2} \mathbf{f}^H \mathbf{A}^{-1} \mathbf{f}, \quad (1)$$

for arbitrary  $\mathbf{z} \in \mathbb{C}^n$ . The right hand side is clearly strictly convex (as a function of  $\mathbf{z}$ ), and has the minimizer  $\mathbf{z} = \mathbf{x} = \mathbf{A}^{-1}\mathbf{f}$ . Since the right and left hand sides must have the same minimizer, we are done.



## Proof of (4), Cont.

(4:  $\Leftarrow$ ): Now suppose that  $\mathbf{x}$  minimizes the function  $E_A$ , and let  $\mathbf{u} \in \mathbb{C}^n$  be an arbitrary unit vector. Now define  $g(s, t) = E_A(\mathbf{x} + \alpha\mathbf{u})$ , where  $\alpha = s + it$ ,  $s, t \in \mathbb{R}$ . Then

$$\begin{aligned}
 g(s, t) &= \frac{1}{2} (\mathbf{x} + \alpha\mathbf{u})^H A (\mathbf{x} + \alpha\mathbf{u}) - \Re \left( (\mathbf{x} + \alpha\mathbf{u})^H \mathbf{f} \right) \\
 &= \frac{1}{2} \mathbf{x}^H A \mathbf{x} + \Re \left( \bar{\alpha} \mathbf{u}^H A \mathbf{x} \right) + \frac{|\alpha|^2}{2} \mathbf{u}^H A \mathbf{u} - \Re \left( \mathbf{x}^H \mathbf{f} \right) - \Re \left( \bar{\alpha} \mathbf{u}^H \mathbf{f} \right) \\
 &= E_A(\mathbf{x}) + \Re \left( \bar{\alpha} \mathbf{u}^H (A \mathbf{x} - \mathbf{f}) \right) + \frac{|\alpha|^2}{2} \mathbf{u}^H A \mathbf{u} \\
 &= E_A(\mathbf{x}) + \Re(\bar{\alpha}) \Re \left( \mathbf{u}^H (A \mathbf{x} - \mathbf{f}) \right) - \Im(\bar{\alpha}) \Im \left( \mathbf{u}^H (A \mathbf{x} - \mathbf{f}) \right) \\
 &\quad + \frac{s^2 + t^2}{2} \mathbf{u}^H A \mathbf{u} \\
 &= E_A(\mathbf{x}) + s \Re \left( \mathbf{u}^H (A \mathbf{x} - \mathbf{f}) \right) + t \Im \left( \mathbf{u}^H (A \mathbf{x} - \mathbf{f}) \right) + \frac{s^2 + t^2}{2} \mathbf{u}^H A \mathbf{u}.
 \end{aligned}$$



### Proof of (4), Cont.

Clearly  $g$  is a strictly convex, quadratic function on  $\mathbb{R}^2$ . Moreover,  $g$  is minimized at  $(s, t) = (0, 0)$ . Hence,

$$0 = \frac{\partial g}{\partial s}(0, 0) = \Re \left( \mathbf{u}^H (\mathbf{A}\mathbf{x} - \mathbf{f}) \right)$$

and

$$0 = \frac{\partial g}{\partial t}(0, 0) = \Im \left( \mathbf{u}^H (\mathbf{A}\mathbf{x} - \mathbf{f}) \right)$$

hold for any vector  $\mathbf{u}$ . It follows then that  $\mathbf{A}\mathbf{x} = \mathbf{f}$ .







## Proposition (inner products and HPD matrices)

*Suppose that  $(\cdot, \cdot) : \mathbb{C}^n \times \mathbb{C}^n \rightarrow \mathbb{C}$  is an inner product. There exists a unique HPD matrix  $A \in \mathbb{C}^{n \times n}$  such that*

$$(\mathbf{x}, \mathbf{y}) = (A\mathbf{x}, \mathbf{y})_2 = (\mathbf{x}, \mathbf{y})_A, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{C}^n.$$

## Proof.

Let, for  $j = 1, \dots, n$ ,  $\hat{\mathbf{e}}_j$  denote the canonical unit basis vectors of  $\mathbb{C}^n$ . Define the matrix  $A = [a_{ij}] \in \mathbb{C}^{n \times n}$  via

$$a_{ij} = (\hat{\mathbf{e}}_j, \hat{\mathbf{e}}_i).$$

The reader should show that  $A$  has the desired properties. □



## Definition (A-conjugate)

Suppose that  $(\cdot, \cdot) : \mathbb{C}^n \times \mathbb{C}^n \rightarrow \mathbb{C}$  is an inner product and  $A \in \mathbb{C}^{n \times n}$  is its associated HPD matrix. We say that  $B \in \mathbb{C}^{n \times n}$  is **self-adjoint** with respect to this inner product iff

$$(\mathbf{x}, B\mathbf{y}) = (\mathbf{x}, B\mathbf{y})_A = (B\mathbf{x}, \mathbf{y})_A = (B\mathbf{x}, \mathbf{y}), \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{C}^n.$$

We say that  $B$  is **self-adjoint positive definite** with respect to the inner product iff  $B$  is self-adjoint and satisfies

$$(\mathbf{x}, B\mathbf{x}) = (\mathbf{x}, B\mathbf{x})_A > 0, \quad \forall \mathbf{x} \in \mathbb{C}_*^n.$$

We say that a set  $S \subset \mathbb{C}^n$  of non-zero vectors is called **A-orthogonal** (or **A-conjugate**) iff whenever  $\mathbf{x}, \mathbf{y} \in S$ , and  $\mathbf{x} \neq \mathbf{y}$ , then

$$(\mathbf{x}, \mathbf{y}) = (\mathbf{x}, \mathbf{y})_A = 0.$$

We say that  $S \subset \mathbb{C}^n$  is **A-orthonormal** iff  $S$  is A-orthogonal and

$$\|\mathbf{x}\|_A = 1, \quad \forall \mathbf{x} \in S.$$



## Proposition (square root)

*Suppose that  $A \in \mathbb{C}^{n \times n}$  is HPD. Then  $A$  is invertible and  $A^{-1}$  is HPD. Furthermore, there exists a unique HPD matrix  $B \in \mathbb{C}^{n \times n}$  with the property that  $BB = A$ .*

## Proof.

Since  $A$  is HPD, it is unitarily diagonalizable, that is, there are a unitary matrix  $U \in \mathbb{C}^{n \times n}$  and a diagonal matrix with positive diagonal entries  $D = \text{diag}[\lambda_1, \dots, \lambda_n]$  such that

$$A = UDU^H.$$

Set

$$B = UD^{1/2}U^H \quad \text{and} \quad D^{1/2} = \text{diag} \left( \sqrt{\lambda_1}, \dots, \sqrt{\lambda_n} \right).$$

Then  $B$  has the desired properties. Furthermore, it is easy to see that

$$A^{-1} = UD^{-1}U^H, \quad D^{-1} = \text{diag} \left( \lambda_1^{-1}, \dots, \lambda_n^{-1} \right).$$

We leave it to the reader to check the details. □



## Proposition (product of HPD matrices)

*Suppose that  $A, B \in \mathbb{C}^{n \times n}$  are HPD matrices. Then the product  $BA$  is self-adjoint and positive definite with respect to  $(\cdot, \cdot)_A$ , i.e.,*

$$(\mathbf{x}, BA\mathbf{y})_A = (BA\mathbf{x}, \mathbf{y})_A, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{C}^n$$

*and*

$$(BA\mathbf{x}, \mathbf{x})_A > 0, \quad \forall \mathbf{x} \in \mathbb{C}_*^n.$$



## Proposition (similarity)

Suppose that  $A \in \mathbb{C}^{n \times n}$  is HPD and  $B \in \mathbb{C}^{n \times n}$  is self-adjoint with respect to  $(\cdot, \cdot)_A$ , i.e.,

$$(\mathbf{x}, B\mathbf{y})_A = (B\mathbf{x}, \mathbf{y})_A, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{C}^n.$$

Then,  $B$  is similar to a matrix that is Hermitian, that is, a matrix that is self-adjoint with respect to the Euclidean inner product  $(\cdot, \cdot)_2$ .

## Proof.

Since  $B \in \mathbb{C}^{n \times n}$  is self-adjoint with respect to  $(\cdot, \cdot)_A$ , then the reader should confirm that

$$B^H A = A B.$$

Since  $A$  is HPD, there is an invertible matrix  $L \in \mathbb{C}^{n \times n}$  such that  $A = LL^H$ . Define

$$C = L^H B L^{-H}.$$



## Proof, Cont.

Then

$$\begin{aligned}C^H &= L^{-1}B^HL \\&= L^{-1}B^HLL^HL^{-H} \\&= L^{-1}B^HAL^{-H} \\&= L^{-1}ABL^{-H} \\&= L^{-1}LL^HBL^{-H} \\&= L^HBL^{-H} \\&= C.\end{aligned}$$

Since  $C$  is similar to  $B$ , the result follows. □



## Theorem (spectral decomposition)

*Suppose that  $A \in \mathbb{C}^{n \times n}$  is HPD and  $B \in \mathbb{C}^{n \times n}$  is self-adjoint with respect to  $(\cdot, \cdot)_A$ . Then all of the eigenvalues of  $B$  are real and there is an  $A$ -orthonormal basis of  $\mathbb{C}^n$  consisting of eigenvectors of  $B$ .*

## Proof.

Applying the Euclidean spectral decomposition theorem to the Hermitian matrix

$$C = L^H B L^{-H},$$

where the invertible matrix  $L$  is taken from the Cholesky-type decomposition  $A = LL^H$  of  $A$ , there are a unitary matrix  $U \in \mathbb{C}^{n \times n}$  and a diagonal matrix with real entries  $D = \text{diag}(\lambda_1, \dots, \lambda_n)$  such that

$$L^H B L^{-H} = C = U D U^{-1}.$$

Hence,  $B$  is similar to a diagonal matrix with real entries:

$$B = (L^{-H} U) D (L^{-H} U)^{-1}.$$



## Proof, Cont.

Moreover, setting  $M = L^{-H}U$ , we see that

$$BM = MD,$$

which implies that the columns of the invertible matrix  $M$  are eigenvectors of  $B$ . It only remains to check that the columns of  $M$  form an  $A$ -orthonormal set. This is left to the reader as an exercise. □

## Corollary (eigenvalues)

*Suppose that  $A \in \mathbb{C}^{n \times n}$  is HPD and  $B \in \mathbb{C}^{n \times n}$  is self-adjoint and positive definite with respect to  $(\cdot, \cdot)_A$ . Then the eigenvalues of  $B$  are all real and positive.*





# Gradient Descent Methods



## Definition (gradient descent)

Suppose that  $A \in \mathbb{C}^{n \times n}$  is HPD,  $\mathbf{f} \in \mathbb{C}^n$ , and define the quadratic function  $E_A : \mathbb{C}^n \rightarrow \mathbb{R}$  via

$$E_A(\mathbf{z}) = \frac{1}{2} \mathbf{z}^H A \mathbf{z} - \Re(\mathbf{z}^H \mathbf{f}).$$

A **gradient descent method** is a two-layer iterative scheme to approximate  $\mathbf{x} = A^{-1}\mathbf{f}$ . Starting from an arbitrary initial guess  $\mathbf{x}_0$ , the iterations proceed as

$$\mathbf{x}_k = \mathbf{x}_{k-1} + \alpha_k \mathbf{d}_{k-1}, \quad k = 1, 2, 3, \dots,$$

where  $\mathbf{d}_{k-1} \in \mathbb{C}^n$  is the  $(k-1)$ -st **search direction**, supplied by the algorithm, and  $\alpha_k \in \mathbb{C}$  is the **step size** given by the condition

$$\alpha_k = \underset{\alpha \in \mathbb{C}}{\operatorname{argmin}} E_A(\mathbf{x}_{k-1} + \alpha \mathbf{d}_{k-1}),$$

which is called a **line search**.



## Theorem (gradient descent, real case)

Suppose that  $A \in \mathbb{R}^{n \times n}$  is SPD,  $\mathbf{f} \in \mathbb{R}^n$ , and define the quadratic function  $E_A : \mathbb{R}^n \rightarrow \mathbb{R}$

$$E_A(\mathbf{z}) = \frac{1}{2} \mathbf{z}^T A \mathbf{z} - \mathbf{z}^T \mathbf{f}.$$

Suppose that the search direction  $\mathbf{d}_{k-1} \in \mathbb{R}_*^n$  and previous iterate  $\mathbf{x}_{k-1} \in \mathbb{R}^n$  in a gradient descent method are given, and define  $\mathbf{r}_{k-1} = \mathbf{f} - A\mathbf{x}_{k-1}$ . Then the step size can be computed exactly via the formula

$$\alpha_k = \operatorname{argmin}_{\alpha \in \mathbb{R}} E_A(\mathbf{x}_{k-1} + \alpha \mathbf{d}_{k-1}) = \frac{\mathbf{d}_{k-1}^T \mathbf{r}_{k-1}}{\mathbf{d}_{k-1}^T A \mathbf{d}_{k-1}}.$$

In other words, a gradient descent method is well defined once a non-trivial search direction is specified.

## Proof.

Consider the quadratic

$$g(\alpha) = E_A(\mathbf{x}_{k-1} + \alpha \mathbf{d}_{k-1}) = E_A(\mathbf{x}_{k-1}) - \alpha \mathbf{d}_{k-1}^T \mathbf{r}_{k-1} + \alpha^2 \frac{1}{2} \mathbf{d}_{k-1}^T A \mathbf{d}_{k-1}.$$



## Proof, Cont.

A calculation with the first derivative locates the extremum:

$$0 = g'(\alpha_k) = -\mathbf{d}_{k-1}^T \mathbf{r}_{k-1} + \alpha_k \mathbf{d}_{k-1}^T \mathbf{A} \mathbf{d}_{k-1},$$

which implies

$$\alpha_k = \frac{\mathbf{d}_{k-1}^T \mathbf{r}_{k-1}}{\mathbf{d}_{k-1}^T \mathbf{A} \mathbf{d}_{k-1}}.$$

The second derivative indicates that this is a minimum:

$$g''(\alpha) = \mathbf{d}_{k-1}^T \mathbf{A} \mathbf{d}_{k-1} > 0,$$

provided  $\mathbf{d}_{k-1} \neq \mathbf{0}$ .





## Theorem (gradient descent, complex case)

Suppose that  $A \in \mathbb{C}^{n \times n}$  is HPD,  $\mathbf{f} \in \mathbb{C}^n$ , and define the quadratic function  $E_A : \mathbb{C}^n \rightarrow \mathbb{R}$

$$E_A(\mathbf{z}) = \frac{1}{2} \mathbf{z}^H A \mathbf{z} - \Re(\mathbf{z}^H \mathbf{f}).$$

Suppose that the search direction  $\mathbf{d}_{k-1} \in \mathbb{C}_*^n$  and previous iterate  $\mathbf{x}_{k-1} \in \mathbb{C}^n$  in a gradient descent method are given, and define  $\mathbf{r}_{k-1} = \mathbf{f} - A\mathbf{x}_{k-1}$ . Then the step size can be computed exactly via the formula

$$\alpha_k = \operatorname{argmin}_{\alpha \in \mathbb{C}} E_A(\mathbf{x}_{k-1} + \alpha \mathbf{d}_{k-1}) = \frac{\mathbf{d}_{k-1}^H \mathbf{r}_{k-1}}{\mathbf{d}_{k-1}^H A \mathbf{d}_{k-1}}.$$

## Proof.

Let  $\alpha = s + it$ , where  $s, t \in \mathbb{R}$ . Define the function

$$g(s, t) = E_A(\mathbf{x}_{k-1} + \alpha \mathbf{d}_{k-1}).$$



## Proof, Cont.

Then,

$$\begin{aligned}
 g(s, t) &= E_A(\mathbf{x}_{k-1} + \alpha \mathbf{d}_{k-1}) \\
 &= E_A(\mathbf{x}_{k-1}) - \Re\left(\bar{\alpha} \mathbf{d}_{k-1}^H \mathbf{r}_{k-1}\right) + \frac{|\alpha|^2}{2} \mathbf{d}_{k-1}^H \mathbf{A} \mathbf{d}_{k-1} \\
 &= E_A(\mathbf{x}_{k-1}) - s \Re\left(\mathbf{d}_{k-1}^H \mathbf{r}_{k-1}\right) - t \Im\left(\mathbf{d}_{k-1}^H \mathbf{r}_{k-1}\right) + \frac{s^2 + t^2}{2} \mathbf{d}_{k-1}^H \mathbf{A} \mathbf{d}_{k-1}.
 \end{aligned}$$

This is a strictly convex quadratic function of two variables. Setting the first derivatives equal to zero, we find

$$\begin{aligned}
 0 &= \frac{\partial g}{\partial s}(s_k, t_k) = -\Re\left(\mathbf{d}_{k-1}^H \mathbf{r}_{k-1}\right) + s_k \mathbf{d}_{k-1}^H \mathbf{A} \mathbf{d}_{k-1}, \\
 0 &= \frac{\partial g}{\partial t}(s_k, t_k) = -\Im\left(\mathbf{d}_{k-1}^H \mathbf{r}_{k-1}\right) + t_k \mathbf{d}_{k-1}^H \mathbf{A} \mathbf{d}_{k-1},
 \end{aligned}$$

which implies that

$$\alpha_k = s_k + it_k = \frac{\mathbf{d}_{k-1}^H \mathbf{r}_{k-1}}{\mathbf{d}_{k-1}^H \mathbf{A} \mathbf{d}_{k-1}}.$$





# The Steepest Descent Method



## Definition (steepest descent)

Suppose that  $A \in \mathbb{C}^{n \times n}$  is HPD, and  $\mathbf{f} \in \mathbb{C}^n$ . The **steepest descent method** is a gradient descent method for which the search direction  $\mathbf{d}_{k-1}$  is defined to be the residual, i.e.,

$$\mathbf{d}_{k-1} = \mathbf{r}_{k-1} = \mathbf{f} - A\mathbf{x}_{k-1},$$

so that the step size is precisely

$$\alpha_k = \frac{\mathbf{r}_{k-1}^H \mathbf{r}_{k-1}}{\mathbf{r}_{k-1}^H A \mathbf{r}_{k-1}}.$$

If  $B \in \mathbb{C}^{n \times n}$  is an HPD matrix, the **B-preconditioned steepest descent** method is a gradient descent method with search direction

$$\mathbf{d}_{k-1} = B^{-1} \mathbf{r}_{k-1}, \tag{2}$$

so that the step size is precisely

$$\alpha_k = \frac{\mathbf{r}_{k-1}^H B^{-1} \mathbf{r}_{k-1}}{\mathbf{r}_{k-1}^H B^{-1} A B^{-1} \mathbf{r}_{k-1}}.$$





## What's the Deal with Preconditioning?

The idea with preconditioning is that the *preconditioner*,  $B$ , should be like  $A$ , but easier to invert. In fact, if we choose  $B = A$ , which is not practical, we would converge in a single iteration, because (2) would yield the error vector as the search direction. One way of realizing preconditioned steepest descent, theoretically, is to observe that it is just normal steepest descent applied to solve the equation

$$B^{-1}Ax = B^{-1}f. \quad (3)$$



## Proposition (orthogonality)

*Suppose that  $A \in \mathbb{C}^{n \times n}$  is HPD, and  $\mathbf{f} \in \mathbb{C}^n$ . Suppose that  $\{\mathbf{x}_k\}_{k=1}^{\infty}$  is computed using the steepest descent method with the starting vector  $\mathbf{x}_0$ . Then the sequence of residual vectors  $\{\mathbf{r}_k\}_{k=1}^{\infty}$ ,  $\mathbf{r}_k = \mathbf{f} - A\mathbf{x}_k$ , has the property that*

$$(\mathbf{r}_k, \mathbf{r}_{k+1})_2 = \mathbf{r}_{k+1}^H \mathbf{r}_k = 0,$$

*for  $k = 0, 1, 2, \dots$*

## Proof.

The proof is an exercise. □

## Remark (orthogonality)

*The last proposition shows that the steepest descent method has the property that its next search direction is always orthogonal to the last. We will see that this can lead to some bad outcomes.*



## Theorem (error equation)

Suppose that  $A \in \mathbb{C}^{n \times n}$  is HPD,  $\mathbf{f} \in \mathbb{C}^n$ , and  $\mathbf{x} = A^{-1}\mathbf{f}$ . Suppose that  $\{\mathbf{x}_k\}_{k=1}^{\infty}$  is computed using the steepest descent method with the starting value  $\mathbf{x}_0 \in \mathbb{C}^n$ . Then the error  $\mathbf{e}_k = \mathbf{x} - \mathbf{x}_k$  satisfies

$$\|\mathbf{e}_{k+1}\|_A^2 = \gamma_k \|\mathbf{e}_k\|_A^2,$$

where

$$\gamma_k = 1 - \frac{(\mathbf{r}_k^H \mathbf{r}_k)^2}{(\mathbf{r}_k^H A \mathbf{r}_k)(\mathbf{r}_k^H A^{-1} \mathbf{r}_k)}.$$

## Proof.

Suppose that  $\mathbf{x} = A^{-1}\mathbf{f}$ . Then, for any  $\mathbf{z} \in \mathbb{C}^n$ ,

$$E_A(\mathbf{z}) = E_A(\mathbf{x}) + \frac{1}{2} \|\mathbf{z} - \mathbf{x}\|_A^2, \quad (4)$$

where, as usual,

$$E_A(\mathbf{z}) = \frac{1}{2} \mathbf{z}^H A \mathbf{z} - \Re(\mathbf{z}^H \mathbf{f}).$$



## Proof, Cont.

Since  $\mathbf{r}_k = \mathbf{f} - \mathbf{A}\mathbf{x}_k$ , and, for the standard steepest descent method,

$$\alpha_{k+1} = \frac{\mathbf{r}_k^H \mathbf{r}_k}{\mathbf{r}_k^H \mathbf{A} \mathbf{r}_k},$$

a brief calculation shows that

$$E_A(\mathbf{x}_{k+1}) = E_A(\mathbf{x}_k + \alpha_{k+1} \mathbf{r}_k) = E_A(\mathbf{x}_k) - \frac{1}{2} \frac{(\mathbf{r}_k^H \mathbf{r}_k)^2}{\mathbf{r}_k^H \mathbf{A} \mathbf{r}_k}. \quad (5)$$

Combining equations (4) and (5) we get

$$\|\mathbf{e}_{k+1}\|_A^2 = \|\mathbf{e}_k\|_A^2 - \frac{(\mathbf{r}_k^H \mathbf{r}_k)^2}{\mathbf{r}_k^H \mathbf{A} \mathbf{r}_k}. \quad (6)$$

Since  $\mathbf{r}_k = \mathbf{A}\mathbf{e}_k$ , we have

$$\|\mathbf{e}_k\|_A^2 = \mathbf{r}_k^H \mathbf{A}^{-1} \mathbf{r}_k. \quad (7)$$

Combining (6) and (7), we get the desired result. □



## Lemma (Kantorovich inequality)

Let the matrix  $A \in \mathbb{C}^{n \times n}$  be HPD with spectrum  $\sigma(A) = \{\lambda_i\}_{i=1}^n$ , with  $0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ , and spectral condition number

$$\kappa = \kappa_2(A) = \frac{\lambda_n}{\lambda_1}.$$

Then, for any  $\mathbf{x} \in \mathbb{C}_*^n$ ,

$$\frac{(\mathbf{x}^H A \mathbf{x}) (\mathbf{x}^H A^{-1} \mathbf{x})}{(\mathbf{x}^H \mathbf{x})^2} \leq \frac{1}{4} \left( \sqrt{\kappa} + \sqrt{\kappa^{-1}} \right)^2.$$

## Proof.

See the book.





## Theorem (convergence)

Suppose that  $A \in \mathbb{C}^{n \times n}$  is HPD,  $\mathbf{f} \in \mathbb{C}^n$ , and  $\mathbf{x} = A^{-1}\mathbf{f}$ . Suppose that  $\{\mathbf{x}_k\}_{k=1}^{\infty}$  is computed using the steepest descent method with the starting value  $\mathbf{x}_0 \in \mathbb{C}^n$ . Then the following estimate holds

$$\|\mathbf{e}_k\|_A \leq \left( \frac{\kappa - 1}{\kappa + 1} \right)^k \|\mathbf{e}_0\|_A,$$

where  $\kappa = \kappa_2(A)$ .

## Proof.

Using the Kantorovich inequality and the previous theorem to show that

$$\gamma_k \leq 1 - \frac{4}{(\kappa^{-1/2} + \kappa^{1/2})^2} = \left( \frac{\kappa - 1}{\kappa + 1} \right)^2,$$

which implies the result. □



## Example

Let us show, by means of an example, that the convergence rate for the steepest descent method, presented in the previous result, is sharp. Let  $A = \text{diag}[a, b] \in \mathbb{R}^{2 \times 2}$ ,  $a > b > 0$ ,  $\mathbf{f} = \mathbf{0}$ , and  $\mathbf{x} = \mathbf{0}$ . We will prove the following: if

$$\mathbf{x}_i = c_i [b, sa]^T \in V_1 = \langle [b, sa]^T \rangle,$$

for some  $c_i \in \mathbb{R}_+$ , where  $s = \pm 1$ , then

$$\mathbf{x}_{i+1} = c_{i+1} [b, -sa]^T \in V_2 = \langle [b, -sa]^T \rangle,$$

for some coefficient  $c_{i+1} \in \mathbb{R}_+$ . In general,  $\mathbf{x}_k \in V_1 \cup V_2$ . First  $\mathbf{r}_i = -c_i [ab, sab]^T$  and  $\mathbf{A}\mathbf{r}_i = -c_i [a^2b, sab^2]^T$ . Then

$$\alpha_{i+1} = \frac{\mathbf{r}_i^T \mathbf{r}_i}{\mathbf{r}_i^T \mathbf{A} \mathbf{r}_i} = \frac{2c_i^2 a^2 b^2}{c_i^2 (a^3 b^2 + a^2 b^3)} = \frac{2}{a + b}.$$



## Example (Cont.)

Thus,

$$\begin{aligned}\mathbf{x}_{i+1} &= \mathbf{x}_i + \alpha_i \mathbf{r}_i \\ &= \begin{bmatrix} c_i b - \frac{2c_i}{a+b} ab \\ c_i sa - \frac{2c_i}{a+b} sab \end{bmatrix} \\ &= \frac{c_i}{a+b} \begin{bmatrix} b(a+b) - 2ab \\ sa(a+b) - 2sab \end{bmatrix} \\ &= \frac{c_i}{a+b} \begin{bmatrix} b^2 - ab \\ sa^2 - sab \end{bmatrix} \\ &= \frac{c_i(b-a)}{a+b} \begin{bmatrix} b \\ -sa \end{bmatrix}\end{aligned}$$





## Example (Cont.)

Hence

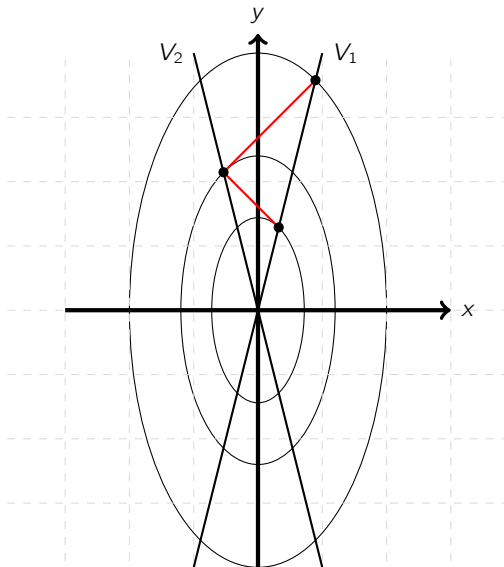
$$c_{i+1} = \frac{c_i(b-a)}{a+b} = -c_i \frac{\kappa-1}{\kappa+1},$$

where  $\kappa = \frac{a}{b}$  is precisely the 2-norm condition number of the matrix  $A$ . From this, one can see that the convergence rate is the worst possible, according to our convergence theory, because

$$c_k = (-1)^k c_0 \left( \frac{\kappa-1}{\kappa+1} \right)^k.$$



Suppose  $a = 16$  and  $b = 4$



## Deterioration of Convergence Rate for Ill-Conditioned Problems



For the steepest descent method, with large spectral condition number  $\kappa$ , we observe that

$$\frac{\kappa - 1}{\kappa + 1} \approx 1 - \frac{2}{\kappa}.$$

In other words, the convergence rate deteriorates as  $\kappa \rightarrow \infty$ , that is, as the system matrix becomes increasingly ill conditioned. We will see, when we study numerical methods for certain types of differential equations, how a family of matrices can become more and more ill conditioned as the size of the matrices increases.



## Theorem (error equation)

Suppose that  $A, B \in \mathbb{C}^{n \times n}$  are HPD,  $\mathbf{f} \in \mathbb{C}^n$ , and  $\mathbf{x} = A^{-1}\mathbf{f}$ . Suppose that  $\{\mathbf{x}_k\}_{k=1}^{\infty}$  is computed using the  $B$ -preconditioned steepest descent method with the starting value  $\mathbf{x}_0 \in \mathbb{C}^n$ . Then the error  $\mathbf{e}_k = \mathbf{x} - \mathbf{x}_k$  satisfies

$$\|\mathbf{e}_{k+1}\|_A^2 = \beta_k \|\mathbf{e}_k\|_A^2,$$

where

$$\beta_k = 1 - \frac{(\mathbf{d}_k^H \mathbf{r}_k)^2}{(\mathbf{d}_k^H A \mathbf{d}_k) (\mathbf{r}_k^H A^{-1} \mathbf{r}_k)} \quad \text{and} \quad \mathbf{r}_k = B \mathbf{d}_k.$$

Setting  $B = L^H L$ , for some invertible matrix  $L \in \mathbb{C}^{n \times n}$ ,

$$\mathbf{g}_k = L \mathbf{d}_k \quad \text{and} \quad C = L^{-H} A L^{-1},$$

the error equation may be expressed using

$$\beta_k = 1 - \frac{(\mathbf{g}_k^H \mathbf{g}_k)^2}{(\mathbf{g}_k^H C \mathbf{g}_k) (\mathbf{g}_k^H C^{-1} \mathbf{g}_k)}.$$



## Theorem (convergence)

Suppose that  $A, B \in \mathbb{C}^{n \times n}$  are HPD,  $\mathbf{f} \in \mathbb{C}^n$ , and  $\mathbf{x} = A^{-1}\mathbf{f}$ . Suppose that  $\{\mathbf{x}_k\}_{k=1}^{\infty}$  is computed using the B-preconditioned steepest descent method with the starting value  $\mathbf{x}_0 \in \mathbb{C}^n$ . Suppose that B has the Cholesky-type factorization  $B = L^H L$ , where  $L \in \mathbb{C}^{n \times n}$  is invertible, and define  $C = L^{-H} A L^{-1}$ . Then the following error estimate holds:

$$\|\mathbf{e}_k\|_A \leq \left( \frac{\kappa_C - 1}{\kappa_C + 1} \right)^k \|\mathbf{e}_0\|_A,$$

where

$$\kappa_C = \kappa_2(C) = \frac{\mu_n}{\mu_1}, \quad \sigma(C) = \{\mu_i\}_{i=1}^n, \quad 0 < \mu_1 \leq \dots \leq \mu_n.$$

## Proof.

Observe that  $C = L^{-H} A L^{-1}$  is HPD. The result now follows by applying the Kantorovich inequality to estimate the size  $\beta_k$  from the previous theorem.  $\square$



## Generalized Condition Number

Consider the preconditioned system  $B^{-1}A\mathbf{x} = B^{-1}\mathbf{f}$ , and observe that  $B^{-1}A$  is self-adjoint and positive definite with respect to  $(\cdot, \cdot)_A$ , since  $B^{-1}$  is HPD. It therefore has positive real eigenvalues. Furthermore, one will find that  $B^{-1}A$  is similar to  $C = L^{-H}AL^{-1}$ . In particular,

$$L(B^{-1}A)L^{-1} = L^{-H}AL^{-1} = C.$$

Therefore the eigenvalues of  $C$  and the preconditioned coefficient matrix  $B^{-1}A$  are the same. Therefore, we will often write the result of the last theorem in the following way:

$$\|\mathbf{e}_k\|_A \leq \left(\frac{\kappa - 1}{\kappa + 1}\right)^k \|\mathbf{e}_0\|_A,$$

where

$$\kappa = \kappa_{B^{-1}A} = \frac{\mu_n}{\mu_1}, \quad \sigma(B^{-1}A) = \{\mu_i\}_{i=1}^n, \quad 0 < \mu_1 \leq \cdots \leq \mu_n.$$

The number  $\kappa_{B^{-1}A} \geq 1$  is called the *condition number* of the preconditioned coefficient matrix  $B^{-1}A$ .

# The Effect of Preconditioning



The result of this last theorem is quite important. It shows that, if we select the preconditioner  $B$  in a careful way, it is possible to dramatically improve the convergence rate. In particular, it is possible, theoretically, to find  $B$  such that  $\kappa_C$  is nearly one. The closer  $\kappa_C$  is to one, the faster is the convergence rate.