

MATH 532: Linear Algebra

Chapter 7: Eigenvalues and Eigenvectors

Greg Fasshauer

Department of Applied Mathematics
Illinois Institute of Technology

Spring 2015



Outline

- 1 Elementary Properties
- 2 Diagonalization via Similarity Transforms
- 3 Functions of Diagonalizable Matrices
- 4 Normal Matrices
- 5 Positive Definite Matrices
- 6 Iterative Solvers
- 7 Krylov Methods



Outline

- 1 Elementary Properties
- 2 Diagonalization via Similarity Transforms
- 3 Functions of Diagonalizable Matrices
- 4 Normal Matrices
- 5 Positive Definite Matrices
- 6 Iterative Solvers
- 7 Krylov Methods



Motivation

Eigenvalues are important, e.g.,

- to decouple systems of ODEs,
- to study physical phenomena such as resonance,
- to tackle the same kind of applications as the SVD (whenever the matrix is symmetric).



Definition

Let A be an $n \times n$ matrix. The scalars λ and nonzero n -vectors \mathbf{x} satisfying

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

are called **eigenvalues** and **eigenvectors** of A . We call (λ, \mathbf{x}) an **eigenpair** of A .



Definition

Let A be an $n \times n$ matrix. The scalars λ and nonzero n -vectors \mathbf{x} satisfying

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

are called **eigenvalues** and **eigenvectors** of A . We call (λ, \mathbf{x}) an **eigenpair** of A .

The set of all eigenvalues of A is called the **spectrum** $\sigma(A)$, i.e.,

$$\sigma(A) = \{\lambda : \lambda \text{ is an eigenvalue of } A\}.$$



Definition

Let A be an $n \times n$ matrix. The scalars λ and nonzero n -vectors \mathbf{x} satisfying

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

are called **eigenvalues** and **eigenvectors** of A . We call (λ, \mathbf{x}) an **eigenpair** of A .

The set of all eigenvalues of A is called the **spectrum** $\sigma(A)$, i.e.,

$$\sigma(A) = \{\lambda : \lambda \text{ is an eigenvalue of } A\}.$$

The **spectral radius** of A is given by

$$\rho(A) = \max_{\lambda \in \sigma(A)} |\lambda|.$$



Theorem

The following are equivalent:

- 1 λ is a eigenvalue of A .
- 2 $A - \lambda I$ is singular.
- 3 $\det(A - \lambda I) = 0$.



Proof.

By definition, λ satisfies $A\mathbf{x} = \lambda\mathbf{x}$. This can be written as

$$(A - \lambda I)\mathbf{x} = \mathbf{0}.$$

We get a nontrivial solution (recall that eigenvectors are always nonzero) if and only if

$A - \lambda I$ is singular.



Proof.

By definition, λ satisfies $A\mathbf{x} = \lambda\mathbf{x}$. This can be written as

$$(A - \lambda I)\mathbf{x} = \mathbf{0}.$$

We get a nontrivial solution (recall that eigenvectors are always nonzero) if and only if

$A - \lambda I$ is singular.

**Remark**

This proof shows that the eigenvector $\mathbf{x} \in N(A - \lambda I)$.



Remark

- In fact, *any* vector in $N(A - \lambda I)$ is an eigenvector of A associated with λ , i.e., eigenvectors are not unique.



Remark

- In fact, *any* vector in $N(A - \lambda I)$ is an eigenvector of A associated with λ , i.e., eigenvectors are not unique.
- Terminology: $N(A - \lambda I)$ is called the *eigenspace* of A associated with λ .



Remark

- In fact, *any* vector in $N(A - \lambda I)$ is an eigenvector of A associated with λ , i.e., eigenvectors are not unique.
- Terminology: $N(A - \lambda I)$ is called the *eigenspace* of A associated with λ .
- *Geometric interpretation: For eigenpairs, matrix multiplication by A acts just like scalar multiplication, i.e., $A\mathbf{x}$ differs from \mathbf{x} only by a stretch factor or a change in orientation (if $\lambda < 0$).*



Definition

Let A be an $n \times n$ matrix. The **characteristic polynomial** of A is given by

$$p(\lambda) = \det(A - \lambda I),$$

and $p(\lambda) = 0$ is called the **characteristic equation**.



Definition

Let A be an $n \times n$ matrix. The **characteristic polynomial** of A is given by

$$p(\lambda) = \det(A - \lambda I),$$

and $p(\lambda) = 0$ is called the **characteristic equation**.

Remark

The basic properties of determinant show that

- *degree(p) = n ,*
- *the leading coefficient, i.e., the coefficient of λ^n is $(-1)^n$.*



Immediate consequences

- 1 The eigenvalues of A are roots of the characteristic polynomial.



Immediate consequences

- 1 The eigenvalues of A are roots of the characteristic polynomial.
- 2 A has n (possibly complex, but necessarily distinct) eigenvalues.



Immediate consequences

- 1 The eigenvalues of A are roots of the characteristic polynomial.
- 2 A has n (possibly complex, but necessarily distinct) eigenvalues.
- 3 If A is real, then complex eigenvalues appear in conjugate pairs, i.e., $\lambda \in \sigma(A) \implies \bar{\lambda} \in \sigma(A)$.



Immediate consequences

- 1 The **eigenvalues** of A are roots of the characteristic polynomial.
- 2 A has n (possibly complex, but necessarily distinct) **eigenvalues**.
- 3 If A is real, then **complex eigenvalues** appear in conjugate pairs, i.e., $\lambda \in \sigma(A) \implies \bar{\lambda} \in \sigma(A)$.
- 4 In particular, simple real (even integer) matrices can have complex eigenvalues and eigenvectors.



Example

Find the eigenvalues and eigenvectors of $A = \begin{pmatrix} 1 & 2 \\ -1 & 1 \end{pmatrix}$.



Example

Find the eigenvalues and eigenvectors of $A = \begin{pmatrix} 1 & 2 \\ -1 & 1 \end{pmatrix}$.

We need to solve

$$p(\lambda) = \det(A - \lambda I) = (1 - \lambda)^2 + 2 = 0$$



Example

Find the eigenvalues and eigenvectors of $A = \begin{pmatrix} 1 & 2 \\ -1 & 1 \end{pmatrix}$.

We need to solve

$$\begin{aligned} p(\lambda) &= \det(A - \lambda I) = (1 - \lambda)^2 + 2 = 0 \\ \iff \lambda^2 - 2\lambda + 3 &= 0 \end{aligned}$$



Example

Find the eigenvalues and eigenvectors of $A = \begin{pmatrix} 1 & 2 \\ -1 & 1 \end{pmatrix}$.

We need to solve

$$\begin{aligned} p(\lambda) &= \det(A - \lambda I) = (1 - \lambda)^2 + 2 = 0 \\ \iff \lambda^2 - 2\lambda + 3 &= 0 \\ \implies \lambda &= \frac{2 \pm \sqrt{4 - 12}}{2} = 1 \pm \sqrt{2}i. \end{aligned}$$



Example

Find the eigenvalues and eigenvectors of $A = \begin{pmatrix} 1 & 2 \\ -1 & 1 \end{pmatrix}$.

We need to solve

$$\begin{aligned} p(\lambda) &= \det(A - \lambda I) = (1 - \lambda)^2 + 2 = 0 \\ \iff \lambda^2 - 2\lambda + 3 &= 0 \\ \implies \lambda &= \frac{2 \pm \sqrt{4 - 12}}{2} = 1 \pm \sqrt{2}i. \end{aligned}$$

Therefore, $\sigma(A) = \{1 + i\sqrt{2}, 1 - i\sqrt{2}\}$.



Example (cont.)

Now, compute the eigenvectors for

$\lambda_1 = 1 + i\sqrt{2}$:

$$A - \lambda_1 I = \begin{pmatrix} -i\sqrt{2} & 2 \\ -1 & -i\sqrt{2} \end{pmatrix} \longrightarrow \begin{pmatrix} 0 & 0 \\ -1 & -i\sqrt{2} \end{pmatrix}$$

so that $N(A - \lambda_1 I) = \text{span}\{(i\sqrt{2}, -1)^T\}$.



Example (cont.)

Now, compute the eigenvectors for

$\lambda_1 = 1 + i\sqrt{2}$:

$$A - \lambda_1 I = \begin{pmatrix} -i\sqrt{2} & 2 \\ -1 & -i\sqrt{2} \end{pmatrix} \longrightarrow \begin{pmatrix} 0 & 0 \\ -1 & -i\sqrt{2} \end{pmatrix}$$

so that $N(A - \lambda_1 I) = \text{span}\{(i\sqrt{2}, -1)^T\}$.

$\lambda_1 = 1 - i\sqrt{2}$:

$$A - \lambda_2 I = \begin{pmatrix} i\sqrt{2} & 2 \\ -1 & i\sqrt{2} \end{pmatrix} \longrightarrow \begin{pmatrix} 0 & 0 \\ -1 & i\sqrt{2} \end{pmatrix}$$

so that $N(A - \lambda_2 I) = \text{span}\{(i\sqrt{2}, 1)^T\}$.



Remark

*Since eigenvalues are the solution of polynomial equations and we know due to **Abel's theorem** that there is **no closed form expression for roots of polynomials of degree five or greater**, **general methods for finding eigenvalues necessarily have to be iterative** (and numerical).*



Formulas for coefficients of characteristic polynomial

If we write

$$(-1)^n p(\lambda) = \lambda^n + c_1 \lambda^{n-1} + c_2 \lambda^{n-2} + \dots + c_{n-1} \lambda + c_n$$

then without proof/derivation (see [Mey00] for details)

$$c_k = (-1)^k s_k, \quad c_0 = 1,$$

where

$$\begin{aligned} s_k &= \sum (\text{all } k \times k \text{ determinant of principal submatrices}) \\ &= \sum (\text{all products of subsets of } k \text{ eigenvalues}) \end{aligned}$$

Special cases

$$\text{trace}(\mathbf{A}) = \lambda_1 + \lambda_2 + \dots + \lambda_n = -c_1,$$

$$\det(\mathbf{A}) = \lambda_1 \lambda_2 \dots \lambda_n = (-1)^n c_n.$$



Example

Compute the characteristic polynomial for

$$A = \begin{pmatrix} 1 & 2 & 1 \\ 0 & -1 & 1 \\ 0 & 0 & 1 \end{pmatrix}$$

We first compute

$$\begin{aligned} (-1)^3 p(\lambda) &= -\det(A - \lambda I) = (1 - \lambda)^2(1 + \lambda) \\ &= (\lambda^2 - 2\lambda + 1)(1 + \lambda) \\ &= \lambda^3 - \lambda^2 - \lambda + 1. \end{aligned}$$



Example (cont.)

On the other hand (using the above formulas)

$$c_0 = 1,$$

$$s_1 = \det(1) = 1 \implies c_1 = -s_1 = -1,$$

$$s_2 = \det \begin{pmatrix} 1 & 2 \\ 0 & -1 \end{pmatrix} + \det \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} + \det \begin{pmatrix} -1 & 1 \\ 0 & 1 \end{pmatrix}$$

$$= -1 + 1 - 1 = -1 \implies c_2 = s_2 = -1,$$

$$s_3 = \det(A) = -1 \implies c_3 = -s_3 = 1.$$



Example (cont.)

The corresponding eigenvectors are

$$\lambda = -1: \mathbf{x} = (1, -1, 0)^T$$

$$\lambda = 1: \mathbf{x} = (1, 0, 0)^T$$

Note that $\lambda = 1$ is a double eigenvalue, but the eigenspace is only one-dimensional, i.e., there is a deficiency (see algebraic vs. geometric multiplicities later).



Example

The trace and determinant combination is particularly applicable to 2×2 problems. Consider

$$A = \begin{pmatrix} 1 & 2 \\ -1 & 1 \end{pmatrix}$$

then

$$\text{trace}(A) = 2 = \lambda_1 + \lambda_2$$

$$\det(A) = 3 = \lambda_1 \lambda_2$$

so that $\lambda_1 = 2 - \lambda_2$ implies

$$(2 - \lambda_2)\lambda_2 = 3 \implies \lambda_2^2 - 2\lambda_2 + 3 = 0$$

as earlier.

Often, the largest eigenvalue is especially important.

Recall spectral radius: $\rho(\mathbf{A}) = \max_{\lambda \in \sigma(\mathbf{A})} |\lambda|$.

A simple upper bound is, using any matrix norm,

$$\rho(\mathbf{A}) \leq \|\mathbf{A}\|.$$

We now prove this.



Proof.

First, we remember submultiplicativity of matrix norms, i.e.,

$$\|AX\| \leq \|A\|\|X\| \quad \text{for any } X. \quad (1)$$

Now, take $X = (\mathbf{x} \ \mathbf{0} \ \dots \ \mathbf{0})$ with (λ, \mathbf{x}) and eigenpair of A . Then $AX = \lambda X$ and

$$\|AX\| = \|\lambda X\| = |\lambda|\|X\|. \quad (2)$$

Combine (1) and (2):

$$\begin{aligned} |\lambda|\|X\| &= \|AX\| \leq \|A\|\|X\| \\ \xRightarrow{\|X\| \neq 0} |\lambda| &\leq \|A\| \\ \xRightarrow{\lambda \text{ arb.}} \rho(A) &\leq \|A\| \end{aligned}$$



More precise estimates of eigenvalues can be obtained with **Gerschgorin circles**.

Definition

Let $A \in \mathbb{C}^{n \times n}$. The **Gerschgorin circles** \mathcal{G}_i of A are defined by

$$\mathcal{G}_i = \{z \in \mathbb{C} : |z - a_{ii}| \leq r_i\}, \quad i = 1, \dots, n$$

with $r_i = \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|$, the (off-diagonal) row sums of A .

Remark

Analogous (but not the same) circles can be defined via column sums.



Theorem

Let $A \in \mathbb{C}^{n \times n}$ and \mathcal{G}_i , $i = 1, \dots, n$, be its Gerschgorin circles. Then

$$\sigma(A) \subseteq \bigcup_{i=1}^n \mathcal{G}_i.$$

Remark

If we use two sets of Gerschgorin circles, \mathcal{G}_r and \mathcal{G}_c (defined via rows sums and via column sums, respectively), then we get a better estimate:

$$\sigma(A) \subseteq \mathcal{G}_r \cap \mathcal{G}_c.$$



Before we prove the theorem we illustrate with an example.

Example

Consider

$$A = \begin{pmatrix} 1 & 0 & 1 \\ 2 & -1 & 0 \\ 1 & 0 & 1 \end{pmatrix}$$

with rough estimate $\rho(A) \leq \|A\|_\infty = 3$.

The Gerschgorin circles are

$$\mathcal{G}_1 = \{z : |z - 1| \leq 1\}$$

$$\mathcal{G}_2 = \{z : |z + 1| \leq 2\}$$

$$\mathcal{G}_3 = \{z : |z - 1| \leq 1\}$$



Proof

Assume (λ, \mathbf{x}) is an eigenpair with \mathbf{x} normalized, i.e., $\|\mathbf{x}\|_\infty = 1$. Consider i such that $|x_i| = \|\mathbf{x}\|_\infty = 1$. Then

$$\lambda x_i = (\lambda \mathbf{x})_i = (\mathbf{A}\mathbf{x})_i = \sum_{j=1}^n a_{ij}x_j = a_{ii}x_i + \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij}x_j$$

so that

$$(\lambda - a_{ii})x_i = \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij}x_j.$$



Proof (cont.)

Then

$$\begin{aligned}
 |\lambda - a_{ii}| &= |\lambda - a_{ii}| \underbrace{|x_i|}_{=1} = \left| \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} x_j \right| \\
 &\stackrel{\Delta \text{ ineq.}}{\leq} \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| \underbrace{|x_j|}_{\leq \|\mathbf{x}\|_\infty = 1} \\
 &\leq \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| = r_i.
 \end{aligned}$$

Therefore $\lambda \in \mathcal{G}_i$ and each λ will lie in some \mathcal{G}_i , i.e.,

$$\sigma(\mathbf{A}) \subseteq \bigcup_{i=1}^n \mathcal{G}_i.$$



Remark

There is no reason to believe that every Gerschgorin circle contains an eigenvalue.

Example

The eigenvalues of $A = \begin{pmatrix} 0 & 1 \\ 4 & 0 \end{pmatrix}$ are $\lambda_{1,2} = \pm 2$.

But we have

$$\mathcal{G}_1 = \{z : |z| \leq 1\}$$

$$\mathcal{G}_2 = \{z : |z| \leq 4\}$$

and \mathcal{G}_1 does not contain an eigenvalue.



Remark

Recall that a *diagonally dominant* matrix satisfies

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad i = 1, \dots, n.$$

However, then the proof above shows that $\lambda = 0$ cannot be an eigenvalue of a diagonally dominant matrix.

Therefore, *diagonally dominant matrices are nonsingular* (cf. HW).



Outline

- 1 Elementary Properties
- 2 Diagonalization via Similarity Transforms**
- 3 Functions of Diagonalizable Matrices
- 4 Normal Matrices
- 5 Positive Definite Matrices
- 6 Iterative Solvers
- 7 Krylov Methods



Recall: **Equivalence**

$A \sim B$ if and only if there exist P, Q nonsingular s.t. $PAQ = B$.



Recall: **Equivalence**

$A \sim B$ if and only if there exist P, Q nonsingular s.t. $PAQ = B$.

Now

Definition

Two $n \times n$ matrices A and B are called **similar** if there exists a nonsingular P such that

$$P^{-1}AP = B.$$



Recall: **Equivalence**

$A \sim B$ if and only if there exist P, Q nonsingular s.t. $PAQ = B$.

Now

Definition

Two $n \times n$ matrices A and B are called **similar** if there exists a nonsingular P such that

$$P^{-1}AP = B.$$

Definition

An $n \times n$ matrix A is called **diagonalizable** if A is similar to a diagonal matrix, i.e., if

$$P^{-1}AP = D$$

for some nonsingular matrix P .

Remark

We already know the SVD, i.e.,

$$A = UDV^T \iff U^T AV = D, \quad U, V \text{ unitary,}$$

where D contains the singular values of A .



Remark

We already know the SVD, i.e.,

$$A = UDV^T \iff U^T AV = D, \quad U, V \text{ unitary,}$$

where D contains the singular values of A .

Now we use a *single transformation matrix*, and D will contain the *eigenvalues* of A .



Remark

We already know the SVD, i.e.,

$$A = UDV^T \iff U^T AV = D, \quad U, V \text{ unitary,}$$

where D contains the singular values of A .

Now we use a *single transformation matrix*, and D will contain the *eigenvalues* of A .

However, *every* matrix A has an SVD. Not so now...



Theorem

An $n \times n$ matrix A is diagonalizable if and only if A possesses a **complete** set of eigenvectors (i.e., it has n linearly independent eigenvectors). Moreover,

$$P^{-1}AP = D = \text{diag}(\lambda_1, \dots, \lambda_n)$$

if and only if (λ_j, P_{*j}) , $j = 1, \dots, n$, are eigenpairs of A .



Theorem

An $n \times n$ matrix A is diagonalizable if and only if A possesses a **complete** set of eigenvectors (i.e., it has n linearly independent eigenvectors). Moreover,

$$P^{-1}AP = D = \text{diag}(\lambda_1, \dots, \lambda_n)$$

if and only if (λ_j, P_{*j}) , $j = 1, \dots, n$, are eigenpairs of A .

Remark

If A possesses a complete set of eigenvectors it is called **nondefective** (or **nondeficient**).



Proof.

$$P^{-1}AP = D = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}$$

$$\iff AP = PD$$

$$\iff A \begin{pmatrix} P_{*1} & P_{*2} & \cdots & P_{*n} \end{pmatrix} = \begin{pmatrix} P_{*1} & P_{*2} & \cdots & P_{*n} \end{pmatrix} \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}$$

$$\iff \begin{pmatrix} AP_{*1} & AP_{*2} & \cdots & AP_{*n} \end{pmatrix} = \begin{pmatrix} \lambda_1 P_{*1} & \lambda_2 P_{*2} & \cdots & \lambda_n P_{*n} \end{pmatrix}$$

$$\iff (\lambda_j, P_{*j}) \text{ is an eigenpair of } A$$

Note that P is invertible if and only if the columns of P are linearly independent. □

Example

Consider

$$A = \begin{pmatrix} 1 & 2 & 1 \\ 0 & -1 & 1 \\ 0 & 0 & 1 \end{pmatrix}$$

with

$$\lambda_1 = 1, \quad N(A - I) = \text{span}\left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \right\}$$

and

$$\lambda_2 = -1, \quad N(A + I) = \text{span}\left\{ \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix} \right\}$$

is **not diagonalizable** since the set of eigenvectors is **not complete**.



Example

Consider

$$A = \begin{pmatrix} 1 & 0 & 1 \\ 2 & -1 & 0 \\ 1 & 0 & 1 \end{pmatrix}$$

with characteristic polynomial

$$p(\lambda) = (1 - \lambda)^2(1 + \lambda) + 1 = \lambda^2 - \lambda^2 - 2\lambda = \lambda(\lambda + 1)(\lambda - 2)$$

and spectrum

$$\sigma(A) = \{-1, 0, 2\}.$$



Example (cont.)

Also, $N(A + I)$:

$$\begin{pmatrix} 2 & 0 & 1 \\ 2 & 0 & 0 \\ 1 & 0 & 2 \end{pmatrix} \longrightarrow \begin{pmatrix} 2 & 0 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & \frac{3}{2} \end{pmatrix}$$

so that $N(A + I) = \text{span}\{(0, 1, 0)^T\}$ (first eigenvector).

Example (cont.)

Also, $N(A + I)$:

$$\begin{pmatrix} 2 & 0 & 1 \\ 2 & 0 & 0 \\ 1 & 0 & 2 \end{pmatrix} \longrightarrow \begin{pmatrix} 2 & 0 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & \frac{3}{2} \end{pmatrix}$$

so that $N(A + I) = \text{span}\{(0, 1, 0)^T\}$ (first eigenvector).

$N(A)$:

$$\begin{pmatrix} 1 & 0 & 1 \\ 2 & -1 & 0 \\ 1 & 0 & 1 \end{pmatrix} \longrightarrow \begin{pmatrix} 1 & 0 & 1 \\ 0 & -1 & -2 \\ 0 & 0 & 0 \end{pmatrix}$$

so that $N(A) = \text{span}\{(-1, -2, 1)^T\}$.

Example (cont.)

Also, $N(A + I)$:

$$\begin{pmatrix} 2 & 0 & 1 \\ 2 & 0 & 0 \\ 1 & 0 & 2 \end{pmatrix} \longrightarrow \begin{pmatrix} 2 & 0 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & \frac{3}{2} \end{pmatrix}$$

so that $N(A + I) = \text{span}\{(0, 1, 0)^T\}$ (first eigenvector). $N(A)$:

$$\begin{pmatrix} 1 & 0 & 1 \\ 2 & -1 & 0 \\ 1 & 0 & 1 \end{pmatrix} \longrightarrow \begin{pmatrix} 1 & 0 & 1 \\ 0 & -1 & -2 \\ 0 & 0 & 0 \end{pmatrix}$$

so that $N(A) = \text{span}\{(-1, -2, 1)^T\}$. $N(A - 2I)$:

$$\begin{pmatrix} -1 & 0 & 1 \\ 2 & -3 & 0 \\ 1 & 0 & -11 \end{pmatrix} \longrightarrow \begin{pmatrix} -1 & 0 & 1 \\ 0 & -3 & 2 \\ 0 & 0 & 0 \end{pmatrix}$$

so that $N(A - 2I) = \text{span}\{(1, \frac{2}{3}, 1)^T\}$.

Example (cont.)

Therefore

$$P = \begin{pmatrix} 0 & -1 & 1 \\ 1 & -2 & \frac{2}{3} \\ 0 & 1 & 1 \end{pmatrix}, \quad \text{so that} \quad P^{-1} = \begin{pmatrix} -\frac{4}{3} & 1 & \frac{2}{3} \\ -\frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{2} & 0 & \frac{1}{2} \end{pmatrix}$$

and

$$P^{-1}AP = \begin{pmatrix} -1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 2 \end{pmatrix}.$$



Theorem

If A, B are similar, then $\sigma(A) = \sigma(B)$.



Theorem

If A, B are similar, then $\sigma(A) = \sigma(B)$.

Proof.

We show $\det(A - \lambda I) = \det(B - \lambda I)$, i.e., A and B have the same characteristic polynomials.

Theorem

If A, B are similar, then $\sigma(A) = \sigma(B)$.

Proof.

We show $\det(A - \lambda I) = \det(B - \lambda I)$, i.e., A and B have the same characteristic polynomials.

Since A, B are similar there exists a nonsingular P such that $P^{-1}AP = B$.

Theorem

If A, B are similar, then $\sigma(A) = \sigma(B)$.

Proof.

We show $\det(A - \lambda I) = \det(B - \lambda I)$, i.e., A and B have the same characteristic polynomials.

Since A, B are similar there exists a nonsingular P such that $P^{-1}AP = B$. Now,

$$\det(B - \lambda I) = \det(P^{-1}AP - \lambda I)$$



Theorem

If A, B are similar, then $\sigma(A) = \sigma(B)$.

Proof.

We show $\det(A - \lambda I) = \det(B - \lambda I)$, i.e., A and B have the same characteristic polynomials.

Since A, B are similar there exists a nonsingular P such that $P^{-1}AP = B$. Now,

$$\begin{aligned}\det(B - \lambda I) &= \det(P^{-1}AP - \lambda I) \\ &= \det(P^{-1}AP - \lambda P^{-1}IP)\end{aligned}$$



Theorem

If A, B are similar, then $\sigma(A) = \sigma(B)$.

Proof.

We show $\det(A - \lambda I) = \det(B - \lambda I)$, i.e., A and B have the same characteristic polynomials.

Since A, B are similar there exists a nonsingular P such that $P^{-1}AP = B$. Now,

$$\begin{aligned}\det(B - \lambda I) &= \det(P^{-1}AP - \lambda I) \\ &= \det(P^{-1}AP - \lambda P^{-1}IP) \\ &= \det\left(P^{-1}(A - \lambda I)P\right)\end{aligned}$$



Theorem

If A, B are similar, then $\sigma(A) = \sigma(B)$.

Proof.

We show $\det(A - \lambda I) = \det(B - \lambda I)$, i.e., A and B have the same characteristic polynomials.

Since A, B are similar there exists a nonsingular P such that $P^{-1}AP = B$. Now,

$$\begin{aligned}\det(B - \lambda I) &= \det(P^{-1}AP - \lambda I) \\ &= \det(P^{-1}AP - \lambda P^{-1}IP) \\ &= \det\left(P^{-1}(A - \lambda I)P\right) \\ &= \det(P^{-1}) \det(A - \lambda I) \det(P) = \det(A - \lambda I)\end{aligned}$$



Theorem

If A, B are similar, then $\sigma(A) = \sigma(B)$.

Proof.

We show $\det(A - \lambda I) = \det(B - \lambda I)$, i.e., A and B have the same characteristic polynomials.

Since A, B are similar there exists a nonsingular P such that $P^{-1}AP = B$. Now,

$$\begin{aligned} \det(B - \lambda I) &= \det(P^{-1}AP - \lambda I) \\ &= \det(P^{-1}AP - \lambda P^{-1}IP) \\ &= \det\left(P^{-1}(A - \lambda I)P\right) \\ &= \det(P^{-1}) \det(A - \lambda I) \det(P) = \det(A - \lambda I) \end{aligned}$$

since $\det(P^{-1}) = \frac{1}{\det(P)}$.



Remark

We saw above that *there exist matrices that are not diagonalizable, i.e., are not similar to a diagonal matrix (of its eigenvalues).*



Remark

We saw above that *there exist matrices that are not diagonalizable, i.e., are not similar to a diagonal matrix (of its eigenvalues).*

However, *every square matrix A is similar to a triangular matrix whose diagonal elements are the eigenvalues of A*

→ *Schur factorization (next).*



Theorem (Schur factorization)

For every $n \times n$ matrix A there exists a unitary matrix U (which is *not unique*) and an upper triangular matrix T (which is also *not unique*) such that

$$U^*AU = T,$$

and the diagonal entries of T are the eigenvalues of A .



Proof

By induction. $n = 1$ is easy: $A = a = \lambda$, $U = 1$, $T = \lambda$.

Assume the statement is true for $n - 1$, and show it also holds for n :

Take (λ, \mathbf{x}) , an eigenpair of A with $\|\mathbf{x}\|_2 = 1$ and construct a Householder reflector R whose first column is \mathbf{x} (see Sect. 5.6), i.e.,

$$\mathbf{x} = R\mathbf{e}_1 \quad \begin{array}{c} R^{-1}=R \\ \longleftrightarrow \end{array} \quad R\mathbf{x} = \mathbf{e}_1.$$

Thus

$$R = (\mathbf{x} \quad V)$$

for some V .



Proof (cont.)

Now

$$\begin{aligned}
 R^*AR &\stackrel{R=R^*}{=} RAR = RA(\mathbf{x} \quad V) \\
 &= R(A\mathbf{x} \quad AV) = R(\lambda\mathbf{x} \quad AV) \\
 &= \left(\lambda \underbrace{R\mathbf{x}}_{=\mathbf{e}_1} \quad RAV \right) = (\lambda\mathbf{e}_1 \quad R^*AV) \\
 &= \begin{pmatrix} \lambda & \mathbf{x}^*AV \\ \mathbf{0} & V^*AV \end{pmatrix}
 \end{aligned}$$

By the induction hypothesis V^*AV is similar to an upper triangular matrix, i.e., there exists a unitary Q such that

$$Q^*(V^*AV)Q = \hat{T}.$$



Proof (cont.)

Finally, let $U = R \begin{pmatrix} 1 & \mathbf{0}^* \\ \mathbf{0} & Q \end{pmatrix}$ so that

$$\begin{aligned}
 U^*AU &= \begin{pmatrix} 1 & \mathbf{0}^* \\ \mathbf{0} & Q^* \end{pmatrix} \underbrace{R^*AR}_{\begin{pmatrix} \lambda & \mathbf{x}^*AV \\ \mathbf{0} & V^*AV \end{pmatrix}} \begin{pmatrix} 1 & \mathbf{0}^* \\ \mathbf{0} & Q \end{pmatrix} \\
 &= \begin{pmatrix} 1 & \mathbf{0}^* \\ \mathbf{0} & Q^* \end{pmatrix} \begin{pmatrix} \lambda & \mathbf{x}^*AVQ \\ \mathbf{0} & V^*AVQ \end{pmatrix} \\
 &= \begin{pmatrix} \lambda & \mathbf{x}^*AVQ \\ \mathbf{0} & \underbrace{Q^*V^*AVQ}_{=\hat{T}} \end{pmatrix} \\
 &= T \quad \text{upper triangular}
 \end{aligned}$$

Proof (cont.)

The diagonal entries of T are the eigenvalues of A since

- the similarity transformation preserves eigenvalues, and
- the eigenvalues of a triangular matrix are its diagonal elements.



Theorem (Cayley–Hamilton Theorem)

Let $A \in \mathbb{C}^{n \times n}$ and let $p(\lambda) = 0$ be its characteristic equation. Then

$$p(A) = 0,$$

i.e., every square matrix satisfies its characteristic equation.



Theorem (Cayley–Hamilton Theorem)

Let $A \in \mathbb{C}^{n \times n}$ and let $p(\lambda) = 0$ be its characteristic equation. Then

$$p(A) = 0,$$

i.e., every square matrix satisfies its characteristic equation.

Proof.

There exist many different proofs. One possibility is via the Schur factorization theorem (see [Mey00, Ex. 7.2.2]). □



Multiplicities

Definition

Let $\lambda \in \sigma(\mathbf{A}) = \{\lambda_1, \lambda_2, \dots, \lambda_k\}$.

- 1 The **algebraic multiplicity** of λ , $\text{algmult}_{\mathbf{A}}(\lambda)$, is its multiplicity as a root of the characteristic equation $p(\lambda) = 0$.
- 2 If $\text{algmult}_{\mathbf{A}}(\lambda) = 1$, then λ is called **simple**.
- 3 The **geometric multiplicity** of λ , $\text{geomult}_{\mathbf{A}}(\lambda)$, is $\dim N(\mathbf{A} - \lambda I)$, the dimension of the eigenspace of λ , i.e., the number of linearly independent eigenvectors associated with λ .
- 4 If $\text{algmult}_{\mathbf{A}}(\lambda) = \text{geomult}_{\mathbf{A}}(\lambda)$, then λ is called **semi-simple**.



Example

Consider

$$A = \begin{pmatrix} -1 & -1 & -2 \\ 8 & -11 & -8 \\ -10 & 11 & 7 \end{pmatrix}, \quad B = \begin{pmatrix} 1 & -4 & -4 \\ 8 & -11 & -8 \\ -8 & 8 & 5 \end{pmatrix}$$

with

$$p_A(\lambda) = p_B(\lambda) = \lambda^3 + 5\lambda^2 + 3\lambda - 9 = (\lambda - 1)(\lambda + 3)^2$$

so that the eigenvalues are

$\lambda = 1$: simple,

$\lambda = -3$: with $\text{algmult}_A(-3) = \text{algmult}_B(-3) = 2$.



Example ((cont.))

Eigenvectors for $\lambda = -3$, A :

$$A + 3I = \begin{pmatrix} 2 & -1 & -2 \\ 8 & -8 & -8 \\ -10 & 11 & 10 \end{pmatrix} \longrightarrow \begin{pmatrix} 2 & -1 & -2 \\ 0 & -4 & 0 \\ 0 & 6 & 0 \end{pmatrix}$$

$$\implies N(A + 3I) = \text{span}\left\{ \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \right\}$$

$$\implies 1 = \text{geomult}_A(-3) < \text{algmult}_A(-3) = 2.$$



Example ((cont.))

Eigenvectors for $\lambda = -3$, B :

$$B + 3I = \begin{pmatrix} 4 & -4 & -4 \\ 8 & -8 & -8 \\ -8 & 8 & 8 \end{pmatrix}$$

$$\implies N(B + 3I) = \text{span} \left\{ \begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} \right\}$$

$$\implies \text{geomult}_B(-3) = 2 = \text{algmult}_B(-3).$$



In general we can say

Theorem

Let $A \in \mathbb{C}^{n \times n}$ and $\lambda \in \sigma(A)$. Then

$$\text{geomult}_A(\lambda) \leq \text{algmult}_A(\lambda).$$



In general we can say

Theorem

Let $A \in \mathbb{C}^{n \times n}$ and $\lambda \in \sigma(A)$. Then

$$\text{geomult}_A(\lambda) \leq \text{algmult}_A(\lambda).$$

Proof

Let's assume that $\text{algmult}_A(\lambda) = k$. If we apply the Schur factorization to A we get

$$U^*AU = \begin{pmatrix} T_{11} & T_{12} \\ \mathbf{0} & T_{22} \end{pmatrix},$$

where T_{11} is $k \times k$ upper triangular with $\text{diag}(T_{11}) = (\lambda, \dots, \lambda)$.



Proof (cont.)

Also, $\lambda \notin \text{diag}(T_{22})$ (where T_{22} is also upper triangular).

Thus $\lambda \notin \sigma(T_{22})$ and

$$T_{22} - \lambda I \text{ is nonsingular,}$$

i.e., $\text{rank}(T_{22} - \lambda I) = n - k$.

Now,

$$\text{geomult}_A(\lambda) = \dim N(A - \lambda I) = n - \text{rank}(A - \lambda I).$$

But, using a unitary (and therefore nonsingular) U ,

$$\begin{aligned} \text{rank}(A - \lambda I) &= \text{rank}(U^*(A - \lambda I)U) \\ &= \text{rank} \begin{pmatrix} T_{11} - \lambda I & T_{12} \\ O & T_{22} - \lambda I \end{pmatrix} \\ &\geq \text{rank}(T_{22} - \lambda I) = n - k. \end{aligned}$$

Therefore

$$\text{geomult}_A(\lambda) \leq n - (n - k) = k = \text{algmult}_A(\lambda). \quad \square$$

Diagonalizability

Theorem

A matrix $A \in \mathbb{C}^{n \times n}$ is diagonalizable if and only if

$$\text{geomult}_A(\lambda) = \text{algmult}_A(\lambda) \quad \text{for all } \lambda \in \sigma(A),$$

i.e., if and only if every eigenvalue is semi-simple.



Diagonalizability

Theorem

A matrix $A \in \mathbb{C}^{n \times n}$ is diagonalizable if and only if

$$\text{geomult}_A(\lambda) = \text{algmult}_A(\lambda) \quad \text{for all } \lambda \in \sigma(A),$$

i.e., if and only if every eigenvalue is semi-simple.

Remark

This provides another interpretation for defective matrices, i.e., a matrix is diagonalizable if and only if it is not defective.



Proof

“ \Leftarrow ”: Assume $\text{geomult}_A(\lambda_i) = \text{algmult}_A(\lambda_i) = a_i$ for all i .
 Furthermore, assume we have k distinct eigenvalues, i.e.,

$$\sigma(\mathbf{A}) = \{\lambda_1, \dots, \lambda_k\}.$$

Take \mathcal{B}_i as a basis for $N(\mathbf{A} - \lambda_i I)$, then

$$\mathcal{B} = \bigcup_{i=1}^k \mathcal{B}_i$$

consists of $\sum_{i=1}^k a_i = n$ vectors.

Moreover, \mathcal{B} is linearly independent (see HW), and it forms a complete set of eigenvectors so that \mathbf{A} is diagonalizable.



Proof (cont.)

“ \implies ”: Assume A is diagonalizable with λ such that $\text{algmult}_A(\lambda) = a$.
Then

$$P^{-1}AP = D = \begin{pmatrix} \lambda I_{a \times a} & \mathbf{O} \\ \mathbf{O} & B \end{pmatrix},$$

where P is nonsingular and B is diagonal with $\lambda \notin B$.
As above,

$$\text{geomult}_A(\lambda) = \dim N(A - \lambda I) = n - \text{rank}(A - \lambda I).$$

However,

$$\begin{aligned} \text{rank}(A - \lambda I) &= \text{rank}(P(D - \lambda I)P^{-1}) \\ &= \text{rank} \begin{pmatrix} \mathbf{O} & \mathbf{O} \\ \mathbf{O} & B - \lambda I \end{pmatrix} = n - a. \end{aligned}$$

Together,

$$\text{geomult}_A(\lambda) = n - (n - a) = \text{algmult}_A(\lambda). \quad \square$$

Corollary

If all eigenvalues of A are simple, then A is diagonalizable.



Corollary

If all eigenvalues of A are simple, then A is diagonalizable.

Remark

*The **converse is not true**. Our earlier example showed that B is diagonalizable since $\sigma(B) = \{-3, 1\}$ with*

$$\begin{aligned}\text{geomult}_B(-3) &= \text{algmult}_B(-3) = 2 \\ \text{geomult}_B(1) &= \text{algmult}_B(1) = 1,\end{aligned}$$

but $\lambda = -3$ is a double eigenvalue.



Spectral Theorem

Theorem

A matrix $A \in \mathbb{C}^{n \times n}$ with $\sigma(A) = \{\lambda_1, \dots, \lambda_k\}$ is diagonalizable if and only if there exist *spectral projectors* G_i , $i = 1, \dots, k$ such that we have the *spectral decomposition*

$$A = \lambda_1 G_1 + \lambda_2 G_2 + \dots + \lambda_k G_k,$$

where the G_i satisfy

- 1 $G_1 + G_2 + \dots + G_k = I$,
- 2 $G_i G_j = O$, $i \neq j$,
- 3 G_i is a projector onto $N(A - \lambda_i I)$ along $R(A - \lambda_i I)$.



Proof

We discuss only “ \implies ” for (1) and (2).

Assume A is diagonalizable, i.e., $A = PDP^{-1}$ with

$$P = (X_1 \quad X_2 \quad \cdots \quad X_k),$$

where the columns of X_i form a basis for $N(A - \lambda_i I)$, i.e.,

$$A = (X_1 \quad X_2 \quad \cdots \quad X_k) \begin{pmatrix} \lambda_1 I & & & 0 \\ & \lambda_2 I & & \\ & & \ddots & \\ 0 & & & \lambda_k I \end{pmatrix} \underbrace{\begin{pmatrix} Y_1^T \\ Y_2^T \\ \vdots \\ Y_k^T \end{pmatrix}}_{=P^{-1}}$$

Proof

We discuss only “ \implies ” for (1) and (2).

Assume A is diagonalizable, i.e., $A = PDP^{-1}$ with

$$P = (X_1 \quad X_2 \quad \cdots \quad X_k),$$

where the columns of X_i form a basis for $N(A - \lambda_i I)$, i.e.,

$$\begin{aligned}
 A &= (X_1 \quad X_2 \quad \cdots \quad X_k) \begin{pmatrix} \lambda_1 I & & & 0 \\ & \lambda_2 I & & \\ & & \ddots & \\ 0 & & & \lambda_k I \end{pmatrix} \underbrace{\begin{pmatrix} Y_1^T \\ Y_2^T \\ \vdots \\ Y_k^T \end{pmatrix}}_{=P^{-1}} \\
 &= \lambda_1 \underbrace{X_1 Y_1^T}_{=G_1} + \lambda_2 \underbrace{X_2 Y_2^T}_{=G_2} + \cdots + \lambda_k \underbrace{X_k Y_k^T}_{=G_k}.
 \end{aligned}$$

Proof (cont.)

The identity

$$A = \lambda_1 \mathbf{G}_1 + \lambda_2 \mathbf{G}_2 + \dots + \lambda_k \mathbf{G}_k$$

is the **spectral decomposition of A** .



Proof (cont.)

The identity

$$A = \lambda_1 G_1 + \lambda_2 G_2 + \dots + \lambda_k G_k$$

is the **spectral decomposition of A** .

If $\lambda_1 = \lambda_2 = \dots = \lambda_k = 1$ then

$$PIP^{-1} = I = G_1 + G_2 + \dots + G_k$$

and we have established (1).



Proof (cont.)

Moreover,

$$P^{-1}P = I \iff \begin{pmatrix} Y_1^T X_1 & Y_1^T X_2 & \cdots & Y_1^T X_k \\ Y_2^T X_1 & Y_2^T X_2 & & \\ & & \ddots & \\ Y_k^T X_1 & & \cdots & Y_k^T X_k \end{pmatrix} = I$$

so that $Y_i^T X_j = \begin{cases} 1, & i = j, \\ 0, & i \neq j, \end{cases}$ and

Proof (cont.)

Moreover,

$$P^{-1}P = I \iff \begin{pmatrix} Y_1^T X_1 & Y_1^T X_2 & \cdots & Y_1^T X_k \\ Y_2^T X_1 & Y_2^T X_2 & & \\ & & \ddots & \\ Y_k^T X_1 & & \cdots & Y_k^T X_k \end{pmatrix} = I$$

so that $Y_i^T X_j = \begin{cases} 1, & i = j, \\ 0, & i \neq j, \end{cases}$ and

$$G_i G_j$$

Proof (cont.)

Moreover,

$$P^{-1}P = I \iff \begin{pmatrix} Y_1^T X_1 & Y_1^T X_2 & \cdots & Y_1^T X_k \\ Y_2^T X_1 & Y_2^T X_2 & & \\ & & \ddots & \\ Y_k^T X_1 & & \cdots & Y_k^T X_k \end{pmatrix} = I$$

so that $Y_i^T X_j = \begin{cases} 1, & i = j, \\ 0, & i \neq j, \end{cases}$ and

$$G_i G_j = X_i \underbrace{Y_i^T X_j}_{=0} Y_j^T$$

Proof (cont.)

Moreover,

$$P^{-1}P = I \iff \begin{pmatrix} Y_1^T X_1 & Y_1^T X_2 & \cdots & Y_1^T X_k \\ Y_2^T X_1 & Y_2^T X_2 & & \\ & & \ddots & \\ Y_k^T X_1 & & \cdots & Y_k^T X_k \end{pmatrix} = I$$

so that $Y_i^T X_j = \begin{cases} 1, & i = j, \\ 0, & i \neq j, \end{cases}$ and

$$G_i G_j = X_i \underbrace{Y_i^T X_j}_{=\delta_{ij}} Y_j^T = \begin{cases} X_i Y_j^T, & i = j, \\ 0, & i \neq j. \end{cases}$$

Proof (cont.)

Moreover,

$$P^{-1}P = I \iff \begin{pmatrix} Y_1^T X_1 & Y_1^T X_2 & \cdots & Y_1^T X_k \\ Y_2^T X_1 & Y_2^T X_2 & & \\ & & \ddots & \\ Y_k^T X_1 & & \cdots & Y_k^T X_k \end{pmatrix} = I$$

so that $Y_i^T X_j = \begin{cases} 1, & i = j, \\ 0, & i \neq j, \end{cases}$ and

$$G_i G_j = X_i \underbrace{Y_i^T X_j}_{=\delta_{ij}} Y_j^T = \begin{cases} X_i Y_j^T, & i = j, \\ 0, & i \neq j. \end{cases}$$

Thus $G_i^2 = G_i$ are projectors and we have established (2). \square

Remark

If λ_j is *simple*, then

$$G_j = \frac{\mathbf{x}\mathbf{y}^*}{\mathbf{y}^*\mathbf{x}},$$

where \mathbf{x} , \mathbf{y}^* , respectively, are the *right and left eigenvectors of A associated with λ_j* .



Outline

- 1 Elementary Properties
- 2 Diagonalization via Similarity Transforms
- 3 Functions of Diagonalizable Matrices**
- 4 Normal Matrices
- 5 Positive Definite Matrices
- 6 Iterative Solvers
- 7 Krylov Methods



Functions of Diagonalizable Matrices

We want to give meaning to

$$f(A),$$

where

A: a square $n \times n$ matrix (below also diagonalizable),

f: a continuous function.



Functions of Diagonalizable Matrices

We want to give meaning to

$$f(A),$$

where

A: a square $n \times n$ matrix (below also diagonalizable),

f: a continuous function.

Functions of matrices play an important role, e.g., in solving systems of ODEs.



Functions of Diagonalizable Matrices

We want to give meaning to

$$f(A),$$

where

A: a square $n \times n$ matrix (below also diagonalizable),

f: a continuous function.

Functions of matrices play an important role, e.g., in solving systems of ODEs.

One possible approach is to use **infinite series**, such as

$$e^A = \sum_{k=0}^{\infty} \frac{A^k}{k!}.$$



Functions of Diagonalizable Matrices

We want to give meaning to

$$f(\mathbf{A}),$$

where

\mathbf{A} : a square $n \times n$ matrix (below also diagonalizable),

f : a continuous function.

Functions of matrices play an important role, e.g., in solving systems of ODEs.

One possible approach is to use **infinite series**, such as

$$e^{\mathbf{A}} = \sum_{k=0}^{\infty} \frac{\mathbf{A}^k}{k!}.$$

However, it is **not so easy to compute this series in practice** (see, e.g., [MVL78, MVL03]) **or to analyze the convergence** of such types of series.



If A is **diagonalizable** then the series are easier to analyze:



If A is **diagonalizable** then the series are easier to analyze:

Recall: A diagonalizable means that there exists a nonsingular P such that

$$P^{-1}AP = D = \text{diag}(\lambda_1, \dots, \lambda_n),$$

where the **eigenvalues** $\lambda_1, \dots, \lambda_n$ need not be distinct.



If A is **diagonalizable** then the series are easier to analyze:

Recall: A diagonalizable means that there exists a nonsingular P such that

$$P^{-1}AP = D = \text{diag}(\lambda_1, \dots, \lambda_n),$$

where the **eigenvalues** $\lambda_1, \dots, \lambda_n$ need not be distinct.

Moreover, from HW 11 we know that

$$P^{-1}A^kP = \text{diag}(\lambda_1^k, \dots, \lambda_n^k) = D^k.$$



If A is **diagonalizable** then the series are easier to analyze:

Recall: A diagonalizable means that there exists a nonsingular P such that

$$P^{-1}AP = D = \text{diag}(\lambda_1, \dots, \lambda_n),$$

where the **eigenvalues** $\lambda_1, \dots, \lambda_n$ need not be distinct.

Moreover, from HW 11 we know that

$$P^{-1}A^kP = \text{diag}(\lambda_1^k, \dots, \lambda_n^k) = D^k.$$

With this setup we can **represent $f(A)$ as a power series** in A .



$$f(A) = \sum_{k=0}^{\infty} c_k A^k$$



$$\begin{aligned} f(A) &= \sum_{k=0}^{\infty} c_k A^k \\ &= \sum_{k=0}^{\infty} c_k P D^k P^{-1} \end{aligned}$$



$$\begin{aligned} f(A) &= \sum_{k=0}^{\infty} c_k A^k \\ &= \sum_{k=0}^{\infty} c_k P D^k P^{-1} = P \left(\sum_{k=0}^{\infty} c_k D^k \right) P^{-1} \end{aligned}$$



$$\begin{aligned}f(A) &= \sum_{k=0}^{\infty} c_k A^k \\&= \sum_{k=0}^{\infty} c_k P D^k P^{-1} = P \left(\sum_{k=0}^{\infty} c_k D^k \right) P^{-1} \\&= P \left(\sum_{k=0}^{\infty} c_k \operatorname{diag}(\lambda_1^k, \dots, \lambda_n^k) \right) P^{-1}\end{aligned}$$



$$\begin{aligned}f(A) &= \sum_{k=0}^{\infty} c_k A^k \\&= \sum_{k=0}^{\infty} c_k P D^k P^{-1} = P \left(\sum_{k=0}^{\infty} c_k D^k \right) P^{-1} \\&= P \left(\sum_{k=0}^{\infty} c_k \operatorname{diag}(\lambda_1^k, \dots, \lambda_n^k) \right) P^{-1} \\&= P \operatorname{diag} \left(\sum_{k=0}^{\infty} c_k \lambda_1^k, \dots, \sum_{k=0}^{\infty} c_k \lambda_n^k \right) P^{-1}\end{aligned}$$



$$\begin{aligned}f(A) &= \sum_{k=0}^{\infty} c_k A^k \\&= \sum_{k=0}^{\infty} c_k P D^k P^{-1} = P \left(\sum_{k=0}^{\infty} c_k D^k \right) P^{-1} \\&= P \left(\sum_{k=0}^{\infty} c_k \operatorname{diag}(\lambda_1^k, \dots, \lambda_n^k) \right) P^{-1} \\&= P \operatorname{diag} \left(\sum_{k=0}^{\infty} c_k \lambda_1^k, \dots, \sum_{k=0}^{\infty} c_k \lambda_n^k \right) P^{-1} \\&= P \operatorname{diag} (f(\lambda_1), \dots, f(\lambda_n)) P^{-1}\end{aligned}$$



$$\begin{aligned}f(A) &= \sum_{k=0}^{\infty} c_k A^k \\&= \sum_{k=0}^{\infty} c_k P D^k P^{-1} = P \left(\sum_{k=0}^{\infty} c_k D^k \right) P^{-1} \\&= P \left(\sum_{k=0}^{\infty} c_k \operatorname{diag}(\lambda_1^k, \dots, \lambda_n^k) \right) P^{-1} \\&= P \operatorname{diag} \left(\sum_{k=0}^{\infty} c_k \lambda_1^k, \dots, \sum_{k=0}^{\infty} c_k \lambda_n^k \right) P^{-1} \\&= P \operatorname{diag} (f(\lambda_1), \dots, f(\lambda_n)) P^{-1} \\&= P f(D) P^{-1}\end{aligned}$$



$$\begin{aligned}
 f(A) &= \sum_{k=0}^{\infty} c_k A^k \\
 &= \sum_{k=0}^{\infty} c_k P D^k P^{-1} = P \left(\sum_{k=0}^{\infty} c_k D^k \right) P^{-1} \\
 &= P \left(\sum_{k=0}^{\infty} c_k \operatorname{diag}(\lambda_1^k, \dots, \lambda_n^k) \right) P^{-1} \\
 &= P \operatorname{diag} \left(\sum_{k=0}^{\infty} c_k \lambda_1^k, \dots, \sum_{k=0}^{\infty} c_k \lambda_n^k \right) P^{-1} \\
 &= P \operatorname{diag} (f(\lambda_1), \dots, f(\lambda_n)) P^{-1} \\
 &= P f(D) P^{-1}
 \end{aligned}$$

Note how the **matrix power series** now has become a diagonal matrix of regular (scalar) power series in the eigenvalues of A .



Thus we can now **define** $f(A)$, A diagonalizable, as

$$\begin{aligned} f(A) &= P f(D) P^{-1} \\ &= P \operatorname{diag}(f(\lambda_1), \dots, f(\lambda_n)) P^{-1}. \end{aligned}$$



Thus we can now **define** $f(A)$, A diagonalizable, as

$$\begin{aligned} f(A) &= P f(D) P^{-1} \\ &= P \operatorname{diag}(f(\lambda_1), \dots, f(\lambda_n)) P^{-1}. \end{aligned}$$

The advantage of this approach is that we have **no problems analyzing convergence of the series** (this is now standard calculus).



Thus we can now **define** $f(A)$, A diagonalizable, as

$$\begin{aligned} f(A) &= P f(D) P^{-1} \\ &= P \operatorname{diag}(f(\lambda_1), \dots, f(\lambda_n)) P^{-1}. \end{aligned}$$

The advantage of this approach is that we have **no problems analyzing convergence of the series** (this is now standard calculus).

However, now there is a **potential problem with uniqueness** since P is not unique.



To understand the uniqueness issue we look more carefully and write

$$f(A) = Pf(D)P^{-1}$$



To understand the uniqueness issue we look more carefully and write

$$\begin{aligned} f(\mathbf{A}) &= \mathbf{P}f(\mathbf{D})\mathbf{P}^{-1} \\ &= (\mathbf{X}_1 \quad \cdots \quad \mathbf{X}_n) \begin{pmatrix} f(\lambda_1)\mathbf{I} & & \\ & \cdots & \\ & & f(\lambda_n)\mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{Y}_1^T \\ \vdots \\ \mathbf{Y}_n^T \end{pmatrix} \end{aligned}$$



To understand the uniqueness issue we look more carefully and write

$$\begin{aligned} f(\mathbf{A}) &= \mathbf{P}f(\mathbf{D})\mathbf{P}^{-1} \\ &= (\mathbf{X}_1 \quad \cdots \quad \mathbf{X}_n) \begin{pmatrix} f(\lambda_1)\mathbf{I} & & \\ & \cdots & \\ & & f(\lambda_n)\mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{Y}_1^T \\ \vdots \\ \mathbf{Y}_n^T \end{pmatrix} \\ &= \sum_{i=1}^n f(\lambda_i)\mathbf{X}_i\mathbf{Y}_i^T \end{aligned}$$



To understand the uniqueness issue we look more carefully and write

$$\begin{aligned}
 f(\mathbf{A}) &= \mathbf{P}f(\mathbf{D})\mathbf{P}^{-1} \\
 &= (\mathbf{X}_1 \quad \cdots \quad \mathbf{X}_n) \begin{pmatrix} f(\lambda_1)\mathbf{I} & & \\ & \cdots & \\ & & f(\lambda_n)\mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{Y}_1^T \\ \vdots \\ \mathbf{Y}_n^T \end{pmatrix} \\
 &= \sum_{i=1}^n f(\lambda_i)\mathbf{X}_i\mathbf{Y}_i^T \\
 &= \sum_{i=1}^n f(\lambda_i)\mathbf{G}_i,
 \end{aligned}$$



To understand the uniqueness issue we look more carefully and write

$$\begin{aligned}
 f(\mathbf{A}) &= \mathbf{P}f(\mathbf{D})\mathbf{P}^{-1} \\
 &= (\mathbf{X}_1 \quad \cdots \quad \mathbf{X}_n) \begin{pmatrix} f(\lambda_1)\mathbf{I} & & \\ & \ddots & \\ & & f(\lambda_n)\mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{Y}_1^T \\ \vdots \\ \mathbf{Y}_n^T \end{pmatrix} \\
 &= \sum_{i=1}^n f(\lambda_i)\mathbf{X}_i\mathbf{Y}_i^T \\
 &= \sum_{i=1}^n f(\lambda_i)\mathbf{G}_i,
 \end{aligned}$$

where the **spectral projectors** \mathbf{G}_i are unique.



To understand the uniqueness issue we look more carefully and write

$$\begin{aligned}
 f(A) &= P f(D) P^{-1} \\
 &= (X_1 \quad \cdots \quad X_n) \begin{pmatrix} f(\lambda_1)I & & \\ & \ddots & \\ & & f(\lambda_n)I \end{pmatrix} \begin{pmatrix} Y_1^T \\ \vdots \\ Y_n^T \end{pmatrix} \\
 &= \sum_{i=1}^n f(\lambda_i) X_i Y_i^T \\
 &= \sum_{i=1}^n f(\lambda_i) G_i,
 \end{aligned}$$

where the **spectral projectors** G_i are unique.

Remark

*Note how the **spectral theorem** helps us convert the problem from one with an infinite series to a single finite sum of length n .*

The representation

$$f(\mathbf{A}) = \sum_{i=1}^n f(\lambda_i) \mathbf{G}_i$$

implies that any function of a diagonalizable matrix \mathbf{A} is a polynomial in \mathbf{A} .



The representation

$$f(\mathbf{A}) = \sum_{i=1}^n f(\lambda_i) \mathbf{G}_i$$

implies that any function of a diagonalizable matrix \mathbf{A} is a polynomial in \mathbf{A} .

To see this, we construct $p(\lambda_j) = f(\lambda_j)$, i.e., we construct a Lagrange interpolating polynomial to f at the eigenvalues of \mathbf{A} :



The representation

$$f(\mathbf{A}) = \sum_{i=1}^n f(\lambda_i) \mathbf{G}_i$$

implies that any function of a diagonalizable matrix \mathbf{A} is a polynomial in \mathbf{A} .

To see this, we construct $p(\lambda_i) = f(\lambda_i)$, i.e., we construct a Lagrange interpolating polynomial to f at the eigenvalues of \mathbf{A} :

$$p(z) = \sum_{i=1}^n f(\lambda_i) L_i(z)$$

$$\text{with } L_i(z) = \prod_{\substack{j=1 \\ j \neq i}}^n (z - \lambda_j) / \prod_{\substack{j=1 \\ j \neq i}}^n (\lambda_i - \lambda_j).$$



Thus,

$$f(\mathbf{A}) = \sum_{i=1}^n f(\lambda_i) \mathbf{G}_i$$



Thus,

$$\begin{aligned} f(\mathbf{A}) &= \sum_{i=1}^n f(\lambda_i) \mathbf{G}_i \\ &= \sum_{i=1}^n p(\lambda_i) \mathbf{G}_i = p(\mathbf{A}). \end{aligned}$$



Thus,

$$\begin{aligned} f(\mathbf{A}) &= \sum_{i=1}^n f(\lambda_i) \mathbf{G}_i \\ &= \sum_{i=1}^n p(\lambda_i) \mathbf{G}_i = p(\mathbf{A}). \end{aligned}$$

On the other hand,

$$p(\mathbf{A}) = \sum_{i=1}^n f(\lambda_i) L_i(\mathbf{A})$$



Thus,

$$\begin{aligned} f(\mathbf{A}) &= \sum_{i=1}^n f(\lambda_i) \mathbf{G}_i \\ &= \sum_{i=1}^n p(\lambda_i) \mathbf{G}_i = p(\mathbf{A}). \end{aligned}$$

On the other hand,

$$p(\mathbf{A}) = \sum_{i=1}^n f(\lambda_i) L_i(\mathbf{A})$$

and we see that

$$\mathbf{G}_i = L_i(\mathbf{A}) = \prod_{\substack{j=1 \\ j \neq i}}^n (\mathbf{A} - \lambda_j \mathbf{I}) / \prod_{\substack{j=1 \\ j \neq i}}^n (\lambda_i - \lambda_j).$$



Remark

- In fact, $f(A)$ is a polynomial in A for *any square* A (see HW — uses Cayley–Hamilton theorem).



Remark

- In fact, $f(A)$ is a polynomial in A for *any square A* (see HW — uses *Cayley–Hamilton theorem*).
- Moreover, for general (square) A we can always *define $f(A)$ via an infinite series*. Then one can prove

Theorem

If $f(z) = \sum_{k=0}^{\infty} c_k(z - z_0)^k$ converges for $|z - z_0| < r$ and $|\lambda_i - z_0| < r$ for all $\lambda_i \in \sigma(A)$, then

$$f(A) = \sum_{k=0}^{\infty} c_k(A - z_0 I)^k.$$



The power method to compute the largest eigenvalue of A

Consider a matrix $A \in \mathbb{C}^{n \times n}$ with eigenvalues

$$|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|,$$

i.e., A has a **dominant** (real) eigenvalue.



The power method to compute the largest eigenvalue of A

Consider a matrix $A \in \mathbb{C}^{n \times n}$ with eigenvalues

$$|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|,$$

i.e., A has a **dominant** (real) eigenvalue.

Note that λ_1 is **real** since if it were complex, then we would also have $\overline{\lambda_1}$ with $|\overline{\lambda_1}| = |\lambda_1|$, so not dominant.



The power method to compute the largest eigenvalue of A

Consider a matrix $A \in \mathbb{C}^{n \times n}$ with eigenvalues

$$|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|,$$

i.e., A has a **dominant** (real) eigenvalue.

Note that λ_1 is **real** since if it were complex, then we would also have $\overline{\lambda_1}$ with $|\overline{\lambda_1}| = |\lambda_1|$, so not dominant.

We now describe a **numerical method to find λ_1** and explain how it can be viewed in the framework of this section.



Consider $f(z) = \left(\frac{z}{\lambda_1}\right)^k$. Then

$$f(\mathbf{A}) = \left(\frac{\mathbf{A}}{\lambda_1}\right)^k$$



Consider $f(z) = \left(\frac{z}{\lambda_1}\right)^k$. Then

$$\begin{aligned} f(\mathbf{A}) &= \left(\frac{\mathbf{A}}{\lambda_1}\right)^k \\ &= \sum_{i=1}^n f(\lambda_i) \mathbf{G}_i \end{aligned}$$



Consider $f(z) = \left(\frac{z}{\lambda_1}\right)^k$. Then

$$\begin{aligned} f(\mathbf{A}) &= \left(\frac{\mathbf{A}}{\lambda_1}\right)^k \\ &= \sum_{i=1}^n f(\lambda_i) \mathbf{G}_i \\ &= \sum_{i=1}^n \left(\frac{\lambda_i}{\lambda_1}\right)^k \mathbf{G}_i \end{aligned}$$



Consider $f(z) = \left(\frac{z}{\lambda_1}\right)^k$. Then

$$\begin{aligned}
 f(\mathbf{A}) &= \left(\frac{\mathbf{A}}{\lambda_1}\right)^k \\
 &= \sum_{i=1}^n f(\lambda_i) \mathbf{G}_i \\
 &= \sum_{i=1}^n \left(\frac{\lambda_i}{\lambda_1}\right)^k \mathbf{G}_i \\
 &= \mathbf{G}_1 + \underbrace{\left(\frac{\lambda_2}{\lambda_1}\right)^k}_{\rightarrow 0} \mathbf{G}_2 + \dots + \underbrace{\left(\frac{\lambda_n}{\lambda_1}\right)^k}_{\rightarrow 0} \mathbf{G}_n \rightarrow \mathbf{G}_1 \quad \text{for } k \rightarrow \infty.
 \end{aligned}$$



Consider $f(z) = \left(\frac{z}{\lambda_1}\right)^k$. Then

$$\begin{aligned}
 f(\mathbf{A}) &= \left(\frac{\mathbf{A}}{\lambda_1}\right)^k \\
 &= \sum_{i=1}^n f(\lambda_i) \mathbf{G}_i \\
 &= \sum_{i=1}^n \left(\frac{\lambda_i}{\lambda_1}\right)^k \mathbf{G}_i \\
 &= \underbrace{\mathbf{G}_1 + \left(\frac{\lambda_2}{\lambda_1}\right)^k \mathbf{G}_2 + \dots + \left(\frac{\lambda_n}{\lambda_1}\right)^k \mathbf{G}_n}_{\rightarrow 0} \rightarrow \mathbf{G}_1 \quad \text{for } k \rightarrow \infty.
 \end{aligned}$$

Therefore

$$\left(\frac{\mathbf{A}}{\lambda_1}\right)^k \mathbf{x}_0 \rightarrow \mathbf{G}_1 \mathbf{x}_0 \in N(\mathbf{A} - \lambda_1 \mathbf{I})$$

since \mathbf{G}_1 is a projector onto $N(\mathbf{A} - \lambda_1 \mathbf{I})$.



Thus any initial vector \mathbf{x}_0 such that $G_1 \mathbf{x}_0 \neq \mathbf{0}$ (i.e., $\mathbf{x}_0 \notin R(A - \lambda_1 I)$) will converge to an eigenvector of A associated with λ_1 via the iteration

$$\frac{A^k \mathbf{x}_0}{\lambda_1^k}, \quad k = 1, 2, \dots$$

In fact, $A^k \mathbf{x}_0$ converges to the first eigenvector, as does any scalar multiple.



Thus any initial vector \mathbf{x}_0 such that $G_1 \mathbf{x}_0 \neq \mathbf{0}$ (i.e., $\mathbf{x}_0 \notin R(A - \lambda_1 I)$) will converge to an eigenvector of A associated with λ_1 via the iteration

$$\frac{A^k \mathbf{x}_0}{\lambda_1^k}, \quad k = 1, 2, \dots$$

In fact, $A^k \mathbf{x}_0$ converges to the first eigenvector, as does any scalar multiple.

To find the eigenvalue λ_1 one iterates for $k = 0, 1, 2, \dots$

$$\mathbf{y}^{(k)} = A\mathbf{x}^{(k)}, \quad \nu^{(k)} = \max\text{comp}(\mathbf{y}^{(k)}), \quad \mathbf{x}^{(k+1)} = \frac{\mathbf{y}^{(k)}}{\nu^{(k)}}.$$



Thus any initial vector \mathbf{x}_0 such that $G_1 \mathbf{x}_0 \neq \mathbf{0}$ (i.e., $\mathbf{x}_0 \notin R(A - \lambda_1 I)$) will converge to an eigenvector of A associated with λ_1 via the iteration

$$\frac{A^k \mathbf{x}_0}{\lambda_1^k}, \quad k = 1, 2, \dots$$

In fact, $A^k \mathbf{x}_0$ converges to the first eigenvector, as does any scalar multiple.

To find the eigenvalue λ_1 one iterates for $k = 0, 1, 2, \dots$

$$\mathbf{y}^{(k)} = A\mathbf{x}^{(k)}, \quad \nu^{(k)} = \max\text{comp}(\mathbf{y}^{(k)}), \quad \mathbf{x}^{(k+1)} = \frac{\mathbf{y}^{(k)}}{\nu^{(k)}}.$$

In fact, $\nu^{(k)} \rightarrow \lambda_1$ since

$$\underbrace{A\mathbf{x}^{(k+1)}}_{\rightarrow A\mathbf{x}_1 = \lambda_1 \mathbf{x}_1} = A \frac{\mathbf{y}^{(k)}}{\nu^{(k)}} = \underbrace{A^2 \mathbf{x}^{(k)}}_{\rightarrow A^2 \mathbf{x}_1 = \lambda_1^2 \mathbf{x}_1} / \nu^{(k)}.$$



Remark

More details of the power method — as well as several other methods for finding eigenvalues — are discussed in MATH 577.



Outline

- 1 Elementary Properties
- 2 Diagonalization via Similarity Transforms
- 3 Functions of Diagonalizable Matrices
- 4 Normal Matrices**
- 5 Positive Definite Matrices
- 6 Iterative Solvers
- 7 Krylov Methods



Normal Matrices

Consider an $n \times n$ matrix A . We know that

- A is diagonalizable (in the sense of similarity) if and only if A is nondefective, and
- A is unitarily similar to a triangular matrix (Schur).



Normal Matrices

Consider an $n \times n$ matrix A . We know that

- A is diagonalizable (in the sense of similarity) if and only if A is nondefective, and
- A is unitarily similar to a triangular matrix (Schur).

Question: What are the conditions on A such that it is unitarily diagonalizable?



Definition

A matrix $A \in \mathbb{C}^{n \times n}$ is called **normal** if

$$A^*A = AA^*.$$



Definition

A matrix $A \in \mathbb{C}^{n \times n}$ is called **normal** if

$$A^*A = AA^*.$$

Theorem

The matrix $A \in \mathbb{C}^{n \times n}$ is *unitarily diagonalizable if and only if it is normal.*



Proof (only easy direction).

Assume A is unitarily diagonalizable, i.e., there exists a unitary U such that

$$U^*AU = D \iff A = UDU^*, A^* = U\bar{D}U^*.$$

Proof (only easy direction).

Assume A is unitarily diagonalizable, i.e., there exists a unitary U such that

$$U^*AU = D \iff A = UDU^*, A^* = U\bar{D}U^*.$$

Then

$$A^*A = U\bar{D}\underbrace{U^*U}_{=I}DU^*,$$



Proof (only easy direction).

Assume A is unitarily diagonalizable, i.e., there exists a unitary U such that

$$U^*AU = D \iff A = UDU^*, A^* = U\bar{D}U^*.$$

Then

$$A^*A = U\bar{D}\underbrace{U^*U}_{=I}DU^*,$$

$$AA^* = UD\underbrace{U^*U}_{=I}\bar{D}U^*.$$



Proof (only easy direction).

Assume A is unitarily diagonalizable, i.e., there exists a unitary U such that

$$U^*AU = D \iff A = UDU^*, A^* = U\bar{D}U^*.$$

Then

$$\begin{aligned} A^*A &= U\bar{D}\underbrace{U^*U}_{=I}DU^*, \\ AA^* &= UD\underbrace{U^*U}_{=I}\bar{D}U^*. \end{aligned}$$

Since

$$\bar{D}D = \sum_{i=1}^n |d_i|^2 = D\bar{D}$$

we have $A^*A = AA^*$ and A is normal. □

Remark

- *Normal matrices are unitarily diagonalizable, i.e., they have an associated complete set of orthogonal eigenvectors.*



Remark

- *Normal matrices are unitarily diagonalizable, i.e., they have an associated complete set of orthogonal eigenvectors.*
- *However, not all complete sets of eigenvectors of normal matrices are orthogonal (see HW).*



Theorem

Let A be normal with $\sigma(A) = \{\lambda_1, \dots, \lambda_k\}$. Then

- 1 $R(A) \perp N(A)$.
- 2 *Eigenvectors to distinct eigenvalues are orthogonal, i.e.,*

$$N(A - \lambda_i I) \perp N(A - \lambda_j I), \quad \lambda_i \neq \lambda_j.$$

- 3 *The spectral projectors G_i are orthogonal projectors.*



Proof

① We know

$$\begin{aligned}N(A^*A) &= N(A), & N(AA^*) &= N(A^*), \\R(A)^\perp &= N(A^*).\end{aligned}$$

Proof

① We know

$$N(A^*A) = N(A), \quad N(AA^*) = N(A^*), \\ R(A)^\perp = N(A^*).$$

Since A is normal we know $N(A^*) = N(A)$ and the statement follows.

Proof

- 1 We know

$$\begin{aligned}N(A^*A) &= N(A), & N(AA^*) &= N(A^*), \\ R(A)^\perp &= N(A^*).\end{aligned}$$

Since A is normal we know $N(A^*) = N(A)$ and the statement follows.

- 2 From above we know that $R(A)^\perp = N(A^*) = N(A)$ whenever A is normal.

Proof

- 1 We know

$$N(A^*A) = N(A), \quad N(AA^*) = N(A^*), \\ R(A)^\perp = N(A^*).$$

Since A is normal we know $N(A^*) = N(A)$ and the statement follows.

- 2 From above we know that $R(A)^\perp = N(A^*) = N(A)$ whenever A is normal.

Moreover, $A - \lambda I$ is also normal since

$$(A - \lambda I)^*(A - \lambda I) = A^*A - \lambda A^* - \bar{\lambda}A + |\lambda|^2 I, \\ (A - \lambda I)(A - \lambda I)^* = AA^* - \bar{\lambda}A - \lambda A^* + |\lambda|^2 I.$$

Therefore,

$$N(A - \lambda I) = N((A - \lambda I)^*) = N(A^* - \bar{\lambda}I).$$

Proof (cont.)

We also have

$$\lambda \in \sigma(\mathbf{A}) \iff \bar{\lambda} \in \sigma(\mathbf{A}^*)$$

since

$$\begin{aligned} \det(\mathbf{A} - \lambda \mathbf{I}) = 0 &\iff \overline{\det(\mathbf{A} - \lambda \mathbf{I})} = 0 \\ &\stackrel{\overline{\det(\mathbf{A})} = \det(\mathbf{A}^*)}{\iff} \det((\mathbf{A} - \lambda \mathbf{I})^*) = 0 \\ &\iff \det(\mathbf{A}^* - \bar{\lambda} \mathbf{I}) = 0. \end{aligned}$$



Proof (cont.)

So we can consider two eigenpairs $(\lambda_i, \mathbf{x}_i)$ and $(\lambda_j, \mathbf{x}_j)$ of A .

Proof (cont.)

So we can consider two eigenpairs $(\lambda_i, \mathbf{x}_i)$ and $(\lambda_j, \mathbf{x}_j)$ of A .
Conjugate transposition yields

$$A\mathbf{x}_j = \lambda_j\mathbf{x}_j \iff \mathbf{x}_j^*A^* = \overline{\lambda_j}\mathbf{x}_j^*,$$

Proof (cont.)

So we can consider two eigenpairs $(\lambda_i, \mathbf{x}_i)$ and $(\lambda_j, \mathbf{x}_j)$ of A .
Conjugate transposition yields

$$A\mathbf{x}_j = \lambda_j\mathbf{x}_j \iff \mathbf{x}_j^*A^* = \overline{\lambda_j}\mathbf{x}_j^*,$$

and from above this is equivalent to

$$\mathbf{x}_j^*A = \lambda_j\mathbf{x}_j^*.$$



Proof (cont.)

So we can consider two eigenpairs $(\lambda_i, \mathbf{x}_i)$ and $(\lambda_j, \mathbf{x}_j)$ of A .
Conjugate transposition yields

$$A\mathbf{x}_j = \lambda_j\mathbf{x}_j \iff \mathbf{x}_j^*A^* = \overline{\lambda_j}\mathbf{x}_j^*,$$

and from above this is equivalent to

$$\mathbf{x}_j^*A = \lambda_j\mathbf{x}_j^*.$$

Now we multiply by \mathbf{x}_i

$$\mathbf{x}_j^* \underbrace{A\mathbf{x}_i}_{=\lambda_i\mathbf{x}_i} = \lambda_j\mathbf{x}_j^*\mathbf{x}_i \iff \lambda_i\mathbf{x}_j^*\mathbf{x}_i = \lambda_j\mathbf{x}_j^*\mathbf{x}_i$$



Proof (cont.)

So we can consider two eigenpairs $(\lambda_i, \mathbf{x}_i)$ and $(\lambda_j, \mathbf{x}_j)$ of A .
Conjugate transposition yields

$$A\mathbf{x}_j = \lambda_j\mathbf{x}_j \iff \mathbf{x}_j^*A^* = \overline{\lambda_j}\mathbf{x}_j^*,$$

and from above this is equivalent to

$$\mathbf{x}_j^*A = \lambda_j\mathbf{x}_j^*.$$

Now we multiply by \mathbf{x}_i

$$\mathbf{x}_j^* \underbrace{A\mathbf{x}_i}_{=\lambda_i\mathbf{x}_i} = \lambda_j\mathbf{x}_j^*\mathbf{x}_i \iff \lambda_i\mathbf{x}_j^*\mathbf{x}_i = \lambda_j\mathbf{x}_j^*\mathbf{x}_i$$

$$\begin{array}{l} \lambda_i \neq \lambda_j \\ \iff \\ \mathbf{x}_j^*\mathbf{x}_i = 0. \end{array}$$



Proof (cont.)

- 3 The **spectral theorem** states that the G_i are projectors onto $N(A - \lambda_i I)$ along $R(A - \lambda_i I)$.



Proof (cont.)

- ③ The **spectral theorem** states that the G_j are projectors onto $N(A - \lambda_j I)$ along $R(A - \lambda_j I)$.

Above we showed that

- $A - \lambda_j I$ is normal provided A is normal, and
- $R(A)^\perp = N(A)$ whenever A is normal.



Proof (cont.)

- ③ The **spectral theorem** states that the G_i are projectors onto $N(A - \lambda_i I)$ along $R(A - \lambda_i I)$.

Above we showed that

- $A - \lambda_i I$ is normal provided A is normal, and
- $R(A)^\perp = N(A)$ whenever A is normal.

Therefore

$$R(A - \lambda_i I)^\perp = N(A - \lambda_i I)$$

and G_i are **orthogonal** projectors. \square



Remark

- *Normal matrices include*
 - *real symmetric, Hermitian, skew-symmetric, skew-Hermitian, orthogonal, and unitary matrices.*



Remark

- *Normal matrices include*
 - *real symmetric, Hermitian, skew-symmetric, skew-Hermitian, orthogonal, and unitary matrices.*
- *All eigenvalues of Hermitian (or real symmetric) matrices are real:*



Remark

- *Normal matrices include*
 - *real symmetric, Hermitian, skew-symmetric, skew-Hermitian, orthogonal, and unitary matrices.*
- *All eigenvalues of Hermitian (or real symmetric) matrices are real:*
First,

$$\mathbf{Ax} = \lambda\mathbf{x} \iff \mathbf{x}^*\mathbf{A}^* = \bar{\lambda}\mathbf{x}^*.$$



Remark

- Normal matrices include
 - real symmetric, Hermitian, skew-symmetric, skew-Hermitian, orthogonal, and unitary matrices.
- All eigenvalues of Hermitian (or real symmetric) matrices are real: First,

$$\mathbf{Ax} = \lambda\mathbf{x} \iff \mathbf{x}^*\mathbf{A}^* = \bar{\lambda}\mathbf{x}^*.$$

Multiply by \mathbf{x}^* and \mathbf{x} , respectively:

$$\mathbf{x}^*\mathbf{Ax} = \lambda\mathbf{x}^*\mathbf{x} \iff \mathbf{x}^*\mathbf{A}^*\mathbf{x} = \bar{\lambda}\mathbf{x}^*\mathbf{x}.$$



Remark

- Normal matrices include
 - real symmetric, Hermitian, skew-symmetric, skew-Hermitian, orthogonal, and unitary matrices.
- All eigenvalues of Hermitian (or real symmetric) matrices are real: First,

$$\mathbf{Ax} = \lambda\mathbf{x} \iff \mathbf{x}^*\mathbf{A}^* = \bar{\lambda}\mathbf{x}^*.$$

Multiply by \mathbf{x}^* and \mathbf{x} , respectively:

$$\mathbf{x}^*\mathbf{Ax} = \lambda\mathbf{x}^*\mathbf{x} \iff \mathbf{x}^*\mathbf{A}^*\mathbf{x} = \bar{\lambda}\mathbf{x}^*\mathbf{x}.$$

Then, since $\mathbf{A}^* = \mathbf{A}$,

$$\lambda\mathbf{x}^*\mathbf{x} = \bar{\lambda}\mathbf{x}^*\mathbf{x} \stackrel{\mathbf{x} \neq 0}{\iff} \lambda = \bar{\lambda}.$$



Moreover, one can show

Theorem

A is real symmetric if and only if A is orthogonally diagonalizable, i.e.,

$$P^T A P = D,$$

where P is orthogonal and D is real.



Rayleigh quotient

Definition

Let $A \in \mathbb{C}^{n \times n}$ and $\mathbf{x} \in \mathbb{C}^n$. Then

$$r(\mathbf{x}) = \frac{\mathbf{x}^* A \mathbf{x}}{\mathbf{x}^* \mathbf{x}}$$

is called the **Rayleigh quotient of A associated with \mathbf{x}** .



Rayleigh quotient

Definition

Let $A \in \mathbb{C}^{n \times n}$ and $\mathbf{x} \in \mathbb{C}^n$. Then

$$r(\mathbf{x}) = \frac{\mathbf{x}^* A \mathbf{x}}{\mathbf{x}^* \mathbf{x}}$$

is called the **Rayleigh quotient of A associated with \mathbf{x}** .

Remark

If \mathbf{x} is an eigenvector of A then $r(\mathbf{x}) = \lambda$, the associated eigenvalue, i.e.,

$$A \mathbf{x} = \lambda \mathbf{x} \quad \implies \quad \mathbf{x}^* A \mathbf{x} = \lambda \mathbf{x}^* \mathbf{x} \quad \iff \quad r(\mathbf{x}) = \lambda.$$



Theorem

Let $A \in \mathbb{C}^{n \times n}$ be Hermitian. Then

$$\lambda_{max} = \max_{\mathbf{x} \neq 0} r(\mathbf{x}), \quad \lambda_{min} = \min_{\mathbf{x} \neq 0} r(\mathbf{x}).$$



Theorem

Let $A \in \mathbb{C}^{n \times n}$ be Hermitian. Then

$$\lambda_{max} = \max_{\mathbf{x} \neq 0} r(\mathbf{x}), \quad \lambda_{min} = \min_{\mathbf{x} \neq 0} r(\mathbf{x}).$$

Remark

Since the eigenvalues of a Hermitian matrix are real they can indeed be ordered.



Theorem

Let $A \in \mathbb{C}^{n \times n}$ be Hermitian. Then

$$\lambda_{\max} = \max_{\mathbf{x} \neq 0} r(\mathbf{x}), \quad \lambda_{\min} = \min_{\mathbf{x} \neq 0} r(\mathbf{x}).$$

Remark

Since the eigenvalues of a Hermitian matrix are real they can indeed be ordered.

Proof (Only for the maximum eigenvalue).

First, we consider an **equivalent formulation**:

$$\lambda_{\max} = \max_{\|\mathbf{x}\|_2=1} \mathbf{x}^* A \mathbf{x}.$$



Proof (cont.)

Now, since A is Hermitian, A is normal and therefore unitarily diagonalizable so that

$$\max_{\|\mathbf{x}\|_2=1} \mathbf{x}^* \mathbf{A} \mathbf{x} = \max_{\|\mathbf{x}\|_2=1} \mathbf{x}^* \mathbf{U} \mathbf{D} \mathbf{U}^* \mathbf{x}.$$

Proof (cont.)

Now, since A is Hermitian, A is normal and therefore unitarily diagonalizable so that

$$\max_{\|\mathbf{x}\|_2=1} \mathbf{x}^* \mathbf{A} \mathbf{x} = \max_{\|\mathbf{x}\|_2=1} \mathbf{x}^* \mathbf{U} \mathbf{D} \mathbf{U}^* \mathbf{x}.$$

Let $\mathbf{y} = \mathbf{U}^* \mathbf{x}$. Then

$$\|\mathbf{y}\|_2 = \|\mathbf{U}^* \mathbf{x}\|_2 = \|\mathbf{x}\|_2$$

and

Proof (cont.)

Now, since A is Hermitian, A is normal and therefore unitarily diagonalizable so that

$$\max_{\|\mathbf{x}\|_2=1} \mathbf{x}^* \mathbf{A} \mathbf{x} = \max_{\|\mathbf{x}\|_2=1} \mathbf{x}^* \mathbf{U} \mathbf{D} \mathbf{U}^* \mathbf{x}.$$

Let $\mathbf{y} = \mathbf{U}^* \mathbf{x}$. Then

$$\|\mathbf{y}\|_2 = \|\mathbf{U}^* \mathbf{x}\|_2 = \|\mathbf{x}\|_2$$

and

$$\max_{\|\mathbf{x}\|_2=1} \mathbf{x}^* \mathbf{A} \mathbf{x} = \max_{\|\mathbf{y}\|_2=1} \mathbf{y}^* \mathbf{D} \mathbf{y}$$

Proof (cont.)

Now, since A is Hermitian, A is normal and therefore unitarily diagonalizable so that

$$\max_{\|\mathbf{x}\|_2=1} \mathbf{x}^* \mathbf{A} \mathbf{x} = \max_{\|\mathbf{x}\|_2=1} \mathbf{x}^* \mathbf{U} \mathbf{D} \mathbf{U}^* \mathbf{x}.$$

Let $\mathbf{y} = \mathbf{U}^* \mathbf{x}$. Then

$$\|\mathbf{y}\|_2 = \|\mathbf{U}^* \mathbf{x}\|_2 = \|\mathbf{x}\|_2$$

and

$$\begin{aligned} \max_{\|\mathbf{x}\|_2=1} \mathbf{x}^* \mathbf{A} \mathbf{x} &= \max_{\|\mathbf{y}\|_2=1} \mathbf{y}^* \mathbf{D} \mathbf{y} \\ &= \max_{\|\mathbf{y}\|_2=1} \sum_{i=1}^n \lambda_i |y_i|^2 \end{aligned}$$

Proof (cont.)

Now, since A is Hermitian, A is normal and therefore unitarily diagonalizable so that

$$\max_{\|\mathbf{x}\|_2=1} \mathbf{x}^* \mathbf{A} \mathbf{x} = \max_{\|\mathbf{x}\|_2=1} \mathbf{x}^* \mathbf{U} \mathbf{D} \mathbf{U}^* \mathbf{x}.$$

Let $\mathbf{y} = \mathbf{U}^* \mathbf{x}$. Then

$$\|\mathbf{y}\|_2 = \|\mathbf{U}^* \mathbf{x}\|_2 = \|\mathbf{x}\|_2$$

and

$$\begin{aligned} \max_{\|\mathbf{x}\|_2=1} \mathbf{x}^* \mathbf{A} \mathbf{x} &= \max_{\|\mathbf{y}\|_2=1} \mathbf{y}^* \mathbf{D} \mathbf{y} \\ &= \max_{\|\mathbf{y}\|_2=1} \sum_{i=1}^n \lambda_i |y_i|^2 \\ &\leq \lambda_{\max} \underbrace{\max_{\|\mathbf{y}\|_2=1} \sum_{i=1}^n |y_i|^2}_{=\|\mathbf{y}\|_2^2} \end{aligned}$$

Proof (cont.)

Now, since A is Hermitian, A is normal and therefore unitarily diagonalizable so that

$$\max_{\|\mathbf{x}\|_2=1} \mathbf{x}^* \mathbf{A} \mathbf{x} = \max_{\|\mathbf{x}\|_2=1} \mathbf{x}^* \mathbf{U} \mathbf{D} \mathbf{U}^* \mathbf{x}.$$

Let $\mathbf{y} = \mathbf{U}^* \mathbf{x}$. Then

$$\|\mathbf{y}\|_2 = \|\mathbf{U}^* \mathbf{x}\|_2 = \|\mathbf{x}\|_2$$

and

$$\begin{aligned} \max_{\|\mathbf{x}\|_2=1} \mathbf{x}^* \mathbf{A} \mathbf{x} &= \max_{\|\mathbf{y}\|_2=1} \mathbf{y}^* \mathbf{D} \mathbf{y} \\ &= \max_{\|\mathbf{y}\|_2=1} \sum_{i=1}^n \lambda_i |y_i|^2 \\ &\leq \lambda_{\max} \underbrace{\max_{\|\mathbf{y}\|_2=1} \sum_{i=1}^n |y_i|^2}_{=\|\mathbf{y}\|_2^2} = \lambda_{\max}. \end{aligned}$$

Proof (cont.)

However, the upper bound can be achieved by making \mathbf{x} a normalized eigenvector for λ_{\max} . Then

$$\mathbf{x}^* \mathbf{A} \mathbf{x} = \mathbf{x}^* \lambda_{\max} \mathbf{x} = \lambda_{\max} \underbrace{\|\mathbf{x}\|_2^2}_{=1} = \lambda_{\max}.$$

So the claim is true. \square



As a generalization one can prove

Theorem (Courant–Fischer Theorem)

Let A be an $n \times n$ Hermitian matrix. Its eigenvalues

$\lambda_{\max} = \lambda_1 \geq \lambda_2 \leq \dots \leq \lambda_n = \lambda_{\min}$ are given by

$$\lambda_i = \max_{\substack{\dim \mathcal{V}=i}} \min_{\substack{\mathbf{x} \in \mathcal{V} \\ \|\mathbf{x}\|_2=1}} \mathbf{x}^* A \mathbf{x}$$

or

$$\lambda_i = \min_{\dim \mathcal{V}=n-i+1} \max_{\substack{\mathbf{x} \in \mathcal{V} \\ \|\mathbf{x}\|_2=1}} \mathbf{x}^* A \mathbf{x}.$$



As a generalization one can prove

Theorem (Courant–Fischer Theorem)

Let A be an $n \times n$ Hermitian matrix. Its eigenvalues

$\lambda_{\max} = \lambda_1 \geq \lambda_2 \leq \dots \leq \lambda_n = \lambda_{\min}$ are given by

$$\lambda_i = \max_{\dim \mathcal{V}=i} \min_{\substack{\mathbf{x} \in \mathcal{V} \\ \|\mathbf{x}\|_2=1}} \mathbf{x}^* A \mathbf{x}$$

or

$$\lambda_i = \min_{\dim \mathcal{V}=n-i+1} \max_{\substack{\mathbf{x} \in \mathcal{V} \\ \|\mathbf{x}\|_2=1}} \mathbf{x}^* A \mathbf{x}.$$

Remark

- Here \mathcal{V} is a subspace of \mathbb{C}^n .
- $i = n$ in the max-min characterization leads to $\mathcal{V} = \mathbb{C}^n$ and λ_{\min} .
- $i = 1$ in the min-max characterization leads to $\mathcal{V} = \mathbb{C}^n$ and λ_{\max} .

Remark

Since the singular values of A are the square roots of the eigenvalues of A^*A an *analogous theorem holds for the singular values of A* (see [Mey00, p. 555] for more details).



Remark

Since the singular values of A are the square roots of the eigenvalues of A^*A an *analogous theorem holds for the singular values of A* (see [Mey00, p. 555] for more details).

In particular,

$$\sigma_{\max} = \max_{\|\mathbf{x}\|_2=1} \mathbf{x}^* A^* A \mathbf{x} = \max_{\|\mathbf{x}\|_2=1} \|\mathbf{A}\mathbf{x}\|_2 = \|A\|_2.$$



Outline

- 1 Elementary Properties
- 2 Diagonalization via Similarity Transforms
- 3 Functions of Diagonalizable Matrices
- 4 Normal Matrices
- 5 Positive Definite Matrices**
- 6 Iterative Solvers
- 7 Krylov Methods



Positive Definite Matrices

Earlier we saw that if $A \in \mathbb{R}^{n \times n}$ is symmetric, then

$$P^T A P = D,$$

where P is an orthogonal matrix of eigenvectors and D is a real diagonal matrix of eigenvalues.



Positive Definite Matrices

Earlier we saw that if $A \in \mathbb{R}^{n \times n}$ is symmetric, then

$$P^T A P = D,$$

where P is an orthogonal matrix of eigenvectors and D is a real diagonal matrix of eigenvalues.

Question: What additional properties of A will ensure that its eigenvalues are all positive (nonnegative)?



A necessary condition

Let's assume that $\lambda_i \geq 0$, $i = 1, \dots, n$. Then

$$\begin{aligned} D &= \text{diag}(\lambda_1, \dots, \lambda_n) \\ &= \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n}) \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n}) = D^{1/2} D^{1/2}. \end{aligned}$$



A necessary condition

Let's assume that $\lambda_i \geq 0$, $i = 1, \dots, n$. Then

$$\begin{aligned} D &= \text{diag}(\lambda_1, \dots, \lambda_n) \\ &= \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n}) \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n}) = D^{1/2} D^{1/2}. \end{aligned}$$

So

$$A = P D P^T = P D^{1/2} D^{1/2} P^T =$$



A necessary condition

Let's assume that $\lambda_i \geq 0$, $i = 1, \dots, n$. Then

$$\begin{aligned} D &= \text{diag}(\lambda_1, \dots, \lambda_n) \\ &= \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n}) \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n}) = D^{1/2} D^{1/2}. \end{aligned}$$

So

$$A = P D P^T = P D^{1/2} D^{1/2} P^T = B^T B,$$

where $B = D^{1/2} P^T$.



A necessary condition

Let's assume that $\lambda_i \geq 0$, $i = 1, \dots, n$. Then

$$\begin{aligned} D &= \text{diag}(\lambda_1, \dots, \lambda_n) \\ &= \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n}) \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n}) = D^{1/2} D^{1/2}. \end{aligned}$$

So

$$A = P D P^T = P D^{1/2} D^{1/2} P^T = B^T B,$$

where $B = D^{1/2} P^T$.

Moreover, $\lambda_i > 0$, $i = 1, \dots, n$, implies D is nonsingular, and therefore B is nonsingular.



A necessary condition

Let's assume that $\lambda_i \geq 0$, $i = 1, \dots, n$. Then

$$\begin{aligned} D &= \text{diag}(\lambda_1, \dots, \lambda_n) \\ &= \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n}) \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n}) = D^{1/2} D^{1/2}. \end{aligned}$$

So

$$A = P D P^T = P D^{1/2} D^{1/2} P^T = B^T B,$$

where $B = D^{1/2} P^T$.

Moreover, $\lambda_i > 0$, $i = 1, \dots, n$, implies **D is nonsingular**, and therefore **B is nonsingular**.

The **converse is also true**, i.e., if **B nonsingular**, then $\lambda_i > 0$ (since $D^{1/2} = B P$ and P orthogonal).



A sufficient condition

Having a factorization

$$A = B^T B$$

is also sufficient:



A sufficient condition

Having a factorization

$$A = B^T B$$

is also sufficient:

Assume (λ, \mathbf{x}) is an eigenpair of A . Then the Rayleigh quotient shows

$$\lambda = \frac{\mathbf{x}^T A \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$$



A sufficient condition

Having a factorization

$$A = B^T B$$

is also sufficient:

Assume (λ, \mathbf{x}) is an eigenpair of A . Then the Rayleigh quotient shows

$$\lambda = \frac{\mathbf{x}^T A \mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \frac{\mathbf{x}^T B^T B \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$$



A sufficient condition

Having a factorization

$$A = B^T B$$

is also sufficient:

Assume (λ, \mathbf{x}) is an eigenpair of A . Then the Rayleigh quotient shows

$$\lambda = \frac{\mathbf{x}^T A \mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \frac{\mathbf{x}^T B^T B \mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \frac{\|B \mathbf{x}\|_2^2}{\|\mathbf{x}\|_2^2} \geq 0.$$



A sufficient condition

Having a factorization

$$A = B^T B$$

is also sufficient:

Assume (λ, \mathbf{x}) is an eigenpair of A . Then the Rayleigh quotient shows

$$\lambda = \frac{\mathbf{x}^T A \mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \frac{\mathbf{x}^T B^T B \mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \frac{\|B \mathbf{x}\|_2^2}{\|\mathbf{x}\|_2^2} \geq 0.$$

Moreover, if B is nonsingular, then $N(B) = \{\mathbf{0}\}$ so that $B \mathbf{x} \neq \mathbf{0}$ and $\lambda > 0$.



A sufficient condition

Having a factorization

$$A = B^T B$$

is also sufficient:

Assume (λ, \mathbf{x}) is an eigenpair of A . Then the Rayleigh quotient shows

$$\lambda = \frac{\mathbf{x}^T A \mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \frac{\mathbf{x}^T B^T B \mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \frac{\|B \mathbf{x}\|_2^2}{\|\mathbf{x}\|_2^2} \geq 0.$$

Moreover, if B is nonsingular, then $N(B) = \{\mathbf{0}\}$ so that $B \mathbf{x} \neq \mathbf{0}$ and $\lambda > 0$.

Conversely, if $\lambda > 0$, then $B \mathbf{x} \neq \mathbf{0}$, and — if $\mathbf{x} \neq \mathbf{0}$ — then B is nonsingular.



Remark

On slide #78 of Chapter 3 we defined:

A symmetric matrix A is positive definite if it has an LU decomposition with positive pivots, i.e.,

$$A = LDL^T = R^T R,$$

where $R = D^{1/2}L^T$ is the upper triangular *Cholesky factor* of A .



Remark

On slide #78 of Chapter 3 we defined:

A symmetric matrix A is positive definite if it has an LU decomposition with positive pivots, i.e.,

$$A = LDL^T = R^T R,$$

where $R = D^{1/2}L^T$ is the upper triangular *Cholesky factor* of A .

This agrees with our discussion above.



Theorem

A real symmetric matrix A is *positive definite* if and only if any of the following equivalent conditions hold:

- 1 A has an LU factorization with positive pivots, or A has a Cholesky factorization $A = R^T R$ with upper triangular matrix R with positive diagonal entries.
- 2 All eigenvalues of A are positive.
- 3 $\mathbf{x}^T A \mathbf{x} > 0$ for all nonzero $\mathbf{x} \in \mathbb{R}^n$.



Remark

- *Earlier we used (1) as the definition of positive definiteness. Often positive definiteness is defined via (3).*



Remark

- *Earlier we used (1) as the definition of positive definiteness. Often positive definiteness is defined via (3).*
- *For a Hermitian matrix A we replace the transpose T by conjugate transpose $*$ and “real” by “complex”.*



Remark

- Earlier we used (1) as the definition of positive definiteness. Often positive definiteness is defined via (3).
- For a Hermitian matrix A we replace the transpose T by conjugate transpose $*$ and “real” by “complex”.
- A few more criteria are listed in [Mey00]. In particular, all principal minors of A must be positive. Therefore, *if A has a nonpositive diagonal entry, then it can't be positive definite.*



Finally,

Definition

Let A be a real symmetric matrix. If

$$\mathbf{x}^T A \mathbf{x} \geq 0$$

for all $\mathbf{x} \in \mathbb{R}^n$, then A is called **positive semidefinite**.



Finally,

Definition

Let A be a real symmetric matrix. If

$$\mathbf{x}^T A \mathbf{x} \geq 0$$

for all $\mathbf{x} \in \mathbb{R}^n$, then A is called **positive semidefinite**.

Theorem

A is *positive semidefinite* if and only if *all eigenvalues of A are nonnegative*.



Finally,

Definition

Let A be a real symmetric matrix. If

$$\mathbf{x}^T A \mathbf{x} \geq 0$$

for all $\mathbf{x} \in \mathbb{R}^n$, then A is called **positive semidefinite**.

Theorem

A is *positive semidefinite* if and only if *all eigenvalues of A are nonnegative*.

Remark

A few more criteria are listed in [Mey00].



Positive definite matrices in applications

- **Gram matrix** in interpolation/least squares approximation:

$$A_{ij} = \langle \mathbf{v}_i, \mathbf{v}_j \rangle$$

where $\{\mathbf{v}_1, \dots, \mathbf{v}_n\} \subseteq \mathcal{V}$, \mathcal{V} some inner product space.



Positive definite matrices in applications

- **Gram matrix** in interpolation/least squares approximation:

$$A_{ij} = \langle \mathbf{v}_i, \mathbf{v}_j \rangle$$

where $\{\mathbf{v}_1, \dots, \mathbf{v}_n\} \subseteq \mathcal{V}$, \mathcal{V} some inner product space.

If the \mathbf{v}_j are linearly independent, then A is positive definite;
otherwise positive semidefinite.



Positive definite matrices in applications

- **Gram matrix** in interpolation/least squares approximation:

$$A_{ij} = \langle \mathbf{v}_i, \mathbf{v}_j \rangle$$

where $\{\mathbf{v}_1, \dots, \mathbf{v}_n\} \subseteq \mathcal{V}$, \mathcal{V} some inner product space.

If the \mathbf{v}_j are linearly independent, then A is positive definite; otherwise positive semidefinite.

- If \mathbf{v}_j are the columns of some matrix V , then $A = V^T V$ is the matrix of the **normal equations** $V^T V \mathbf{x} = V^T \mathbf{b}$.



Positive definite matrices in applications

- **Gram matrix** in interpolation/least squares approximation:

$$A_{ij} = \langle \mathbf{v}_i, \mathbf{v}_j \rangle$$

where $\{\mathbf{v}_1, \dots, \mathbf{v}_n\} \subseteq \mathcal{V}$, \mathcal{V} some inner product space.

If the \mathbf{v}_j are linearly independent, then A is positive definite; otherwise positive semidefinite.

- If \mathbf{v}_i are the columns of some matrix V , then $A = V^T V$ is the matrix of the **normal equations** $V^T V \mathbf{x} = V^T \mathbf{b}$.
- If $\mathbf{v}_i = K(\cdot, \mathbf{x}_i)$ is a (reproducing) kernel function centered at \mathbf{x}_i , then $A_{ij} = \langle K(\cdot, \mathbf{x}_i), K(\cdot, \mathbf{x}_j) \rangle_{\mathcal{H}_K} = K(\mathbf{x}_i, \mathbf{x}_j)$. This is the matrix that appears in **kriging** and **RBF interpolation**.



- Hessian matrix in optimization:



- **Hessian matrix** in optimization: Start with n -dimensional Taylor theorem:

$$f(\mathbf{x}) = f(\mathbf{z}) + \sum_{i=1}^n (x_i - z_i) \frac{\partial f}{\partial x_i}(\mathbf{z}) + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (x_i - z_i)(x_j - z_j) \frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{z}) + \dots$$



- **Hessian matrix** in optimization: Start with n -dimensional Taylor theorem:

$$\begin{aligned}
 f(\mathbf{x}) &= f(\mathbf{z}) + \sum_{i=1}^n (x_i - z_i) \frac{\partial f}{\partial x_i}(\mathbf{z}) + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (x_i - z_i)(x_j - z_j) \frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{z}) + \dots \\
 &= f(\mathbf{z}) + (\mathbf{x} - \mathbf{z})^T \nabla_f(\mathbf{z}) + \frac{1}{2} (\mathbf{x} - \mathbf{z})^T \mathbf{H}_f(\mathbf{z}) (\mathbf{x} - \mathbf{z}) + \dots,
 \end{aligned}$$

where ∇_f is the **gradient** of f and \mathbf{H}_f is its **Hessian matrix**.



- **Hessian matrix** in optimization: Start with **n -dimensional Taylor theorem**:

$$\begin{aligned}
 f(\mathbf{x}) &= f(\mathbf{z}) + \sum_{i=1}^n (x_i - z_i) \frac{\partial f}{\partial x_i}(\mathbf{z}) + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (x_i - z_i)(x_j - z_j) \frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{z}) + \dots \\
 &= f(\mathbf{z}) + (\mathbf{x} - \mathbf{z})^T \nabla_f(\mathbf{z}) + \frac{1}{2} (\mathbf{x} - \mathbf{z})^T \mathbf{H}_f(\mathbf{z}) (\mathbf{x} - \mathbf{z}) + \dots,
 \end{aligned}$$

where ∇_f is the **gradient** of f and \mathbf{H}_f is its **Hessian matrix**.

From calculus it is known that **convexity/concavity at a critical point \mathbf{z}** , i.e., $\nabla_f(\mathbf{z}) = \mathbf{0}$, can be **determined by the Hessian matrix**. In fact,



- **Hessian matrix** in optimization: Start with **n -dimensional Taylor theorem**:

$$\begin{aligned}
 f(\mathbf{x}) &= f(\mathbf{z}) + \sum_{i=1}^n (x_i - z_i) \frac{\partial f}{\partial x_i}(\mathbf{z}) + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (x_i - z_i)(x_j - z_j) \frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{z}) + \dots \\
 &= f(\mathbf{z}) + (\mathbf{x} - \mathbf{z})^T \nabla_f(\mathbf{z}) + \frac{1}{2} (\mathbf{x} - \mathbf{z})^T \mathbf{H}_f(\mathbf{z})(\mathbf{x} - \mathbf{z}) + \dots,
 \end{aligned}$$

where ∇_f is the **gradient** of f and \mathbf{H}_f is its **Hessian matrix**.

From calculus it is known that **convexity/concavity at a critical point \mathbf{z}** , i.e., $\nabla_f(\mathbf{z}) = \mathbf{0}$, can be **determined by the Hessian matrix**. In fact,

- If $\mathbf{H}_f(\mathbf{z})$ is **positive definite**, then f has a **minimum at \mathbf{z}** .



- **Hessian matrix** in optimization: Start with **n -dimensional Taylor theorem**:

$$\begin{aligned}
 f(\mathbf{x}) &= f(\mathbf{z}) + \sum_{i=1}^n (x_i - z_i) \frac{\partial f}{\partial x_i}(\mathbf{z}) + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (x_i - z_i)(x_j - z_j) \frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{z}) + \dots \\
 &= f(\mathbf{z}) + (\mathbf{x} - \mathbf{z})^T \nabla_f(\mathbf{z}) + \frac{1}{2} (\mathbf{x} - \mathbf{z})^T \mathbf{H}_f(\mathbf{z}) (\mathbf{x} - \mathbf{z}) + \dots,
 \end{aligned}$$

where ∇_f is the **gradient** of f and \mathbf{H}_f is its **Hessian matrix**.

From calculus it is known that **convexity/concavity at a critical point \mathbf{z}** , i.e., $\nabla_f(\mathbf{z}) = \mathbf{0}$, can be **determined by the Hessian matrix**. In fact,

- If $\mathbf{H}_f(\mathbf{z})$ is **positive definite**, then f has a **minimum** at \mathbf{z} .
- If $\mathbf{H}_f(\mathbf{z})$ is **negative definite**, then f has a **maximum** at \mathbf{z} .



- **Hessian matrix** in optimization: Start with n -dimensional Taylor theorem:

$$\begin{aligned}
 f(\mathbf{x}) &= f(\mathbf{z}) + \sum_{i=1}^n (x_i - z_i) \frac{\partial f}{\partial x_i}(\mathbf{z}) + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (x_i - z_i)(x_j - z_j) \frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{z}) + \dots \\
 &= f(\mathbf{z}) + (\mathbf{x} - \mathbf{z})^T \nabla_f(\mathbf{z}) + \frac{1}{2} (\mathbf{x} - \mathbf{z})^T \mathbf{H}_f(\mathbf{z})(\mathbf{x} - \mathbf{z}) + \dots,
 \end{aligned}$$

where ∇_f is the **gradient** of f and \mathbf{H}_f is its **Hessian matrix**.

From calculus it is known that **convexity/concavity at a critical point** \mathbf{z} , i.e., $\nabla_f(\mathbf{z}) = \mathbf{0}$, can be **determined by the Hessian matrix**. In fact,

- If $\mathbf{H}_f(\mathbf{z})$ is positive definite, then f has a minimum at \mathbf{z} .
- If $\mathbf{H}_f(\mathbf{z})$ is negative definite, then f has a maximum at \mathbf{z} .

Moreover, if $\mathbf{H}_f(\mathbf{z})$ is positive semidefinite for all points in the domain of f , then f is a convex function.



- Covariance matrix in probability/statistics:



- **Covariance matrix** in probability/statistics: Let $\mathbf{X} = (X_1, \dots, X_n)^T$ be a **vector of random variables** with **mean** $\mu_i = \mathbb{E}[X_i]$, $i = 1, \dots, n$. Then the **covariance matrix of \mathbf{X}** is given by

$$A_{ij} = \mathbb{E}[(X_i - \mu_i)(X_j - \mu_j)]$$



- Covariance matrix** in probability/statistics: Let $\mathbf{X} = (X_1, \dots, X_n)^T$ be a **vector of random variables** with **mean** $\mu_i = \mathbb{E}[X_i]$, $i = 1, \dots, n$. Then the **covariance matrix of \mathbf{X}** is given by

$$A_{ij} = \mathbb{E}[(X_i - \mu_i)(X_j - \mu_j)]$$

We can see that **A is positive semidefinite**:

$$\mathbf{z}^T \mathbf{A} \mathbf{z} = \mathbb{E} \left[\sum_{i=1}^n \sum_{j=1}^n z_i (X_i - \mu_i) (X_j - \mu_j) z_j \right]$$



- Covariance matrix** in probability/statistics: Let $\mathbf{X} = (X_1, \dots, X_n)^T$ be a **vector of random variables** with **mean** $\mu_i = \mathbb{E}[X_i]$, $i = 1, \dots, n$. Then the **covariance matrix of \mathbf{X}** is given by

$$A_{ij} = \mathbb{E}[(X_i - \mu_i)(X_j - \mu_j)]$$

We can see that **A is positive semidefinite**:

$$\begin{aligned} \mathbf{z}^T \mathbf{A} \mathbf{z} &= \mathbb{E} \left[\sum_{i=1}^n \sum_{j=1}^n z_i (X_i - \mu_i) (X_j - \mu_j) z_j \right] \\ &= \mathbb{E} \left[\left(\sum_{i=1}^n z_i (X_i - \mu_i) \right)^2 \right] \end{aligned}$$



- Covariance matrix** in probability/statistics: Let $\mathbf{X} = (X_1, \dots, X_n)^T$ be a **vector of random variables** with **mean** $\mu_i = \mathbb{E}[X_i]$, $i = 1, \dots, n$. Then the **covariance matrix of \mathbf{X}** is given by

$$A_{ij} = \mathbb{E}[(X_i - \mu_i)(X_j - \mu_j)]$$

We can see that **A is positive semidefinite**:

$$\begin{aligned} \mathbf{z}^T \mathbf{A} \mathbf{z} &= \mathbb{E} \left[\sum_{i=1}^n \sum_{j=1}^n z_i (X_i - \mu_i) (X_j - \mu_j) z_j \right] \\ &= \mathbb{E} \left[\left(\sum_{i=1}^n z_i (X_i - \mu_i) \right)^2 \right] \geq 0. \end{aligned}$$



- **Finite difference matrices:** See, e.g., [Mey00, Example 7.6.2].



- **Finite difference matrices:** See, e.g., [Mey00, Example 7.6.2].
- **“Stiffness” matrices:** in finite element formulations, based on the interpretation of **energy of some state \mathbf{x} as a quadratic form $\mathbf{x}^T \mathbf{A} \mathbf{x}$** . Positive energy (a fundamental physical assumption) means positive definite \mathbf{A} .

More details in MATH 581.



Quadratic forms

Definition

Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $\mathbf{x} \in \mathbb{R}^n$. The scalar function

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j$$

is called a **quadratic form**.

The quadratic form $\mathbf{x}^T \mathbf{A} \mathbf{x}$ is called **positive definite** if the matrix \mathbf{A} is positive definite.



Remark

We always assume that the matrix of a quadratic form is symmetric:



Remark

We *always assume that the matrix of a quadratic form is symmetric:*

Even if A is not symmetric, $\frac{A+A^T}{2}$ always is symmetric.



Remark

We *always assume that the matrix of a quadratic form is symmetric:*

Even if A is not symmetric, $\frac{A+A^T}{2}$ always is symmetric.

And we have *for the quadratic form*

$$\mathbf{x}^T \left(\frac{A + A^T}{2} \right) \mathbf{x} = \frac{1}{2} \mathbf{x}^T A \mathbf{x} + \frac{1}{2} \mathbf{x}^T A^T \mathbf{x}$$



Remark

We *always assume that the matrix of a quadratic form is symmetric:*

Even if A is not symmetric, $\frac{A+A^T}{2}$ always is symmetric.

And we have *for the quadratic form*

$$\begin{aligned}\mathbf{x}^T \left(\frac{A + A^T}{2} \right) \mathbf{x} &= \frac{1}{2} \mathbf{x}^T A \mathbf{x} + \frac{1}{2} \mathbf{x}^T A^T \mathbf{x} \\ &= \mathbf{x}^T A \mathbf{x}\end{aligned}$$



Remark

We *always assume that the matrix of a quadratic form is symmetric*:

Even if A is not symmetric, $\frac{A+A^T}{2}$ always is symmetric.

And we have *for the quadratic form*

$$\begin{aligned}\mathbf{x}^T \left(\frac{A + A^T}{2} \right) \mathbf{x} &= \frac{1}{2} \mathbf{x}^T A \mathbf{x} + \frac{1}{2} \mathbf{x}^T A^T \mathbf{x} \\ &= \mathbf{x}^T A \mathbf{x}\end{aligned}$$

because $\mathbf{x}^T A^T \mathbf{x} = \mathbf{x}^T A \mathbf{x}$ is a scalar.



Every quadratic form can be written in **standard** (i.e., diagonal) **form** since every real symmetric matrix is orthogonally similar to a diagonal matrix.



Every quadratic form can be written in **standard** (i.e., diagonal) **form** since every real symmetric matrix is orthogonally similar to a diagonal matrix.

Example

Take

$$\begin{aligned} f(\mathbf{x}) &= x_1 x_2 = \mathbf{x}^T \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} \mathbf{x} \\ &= \mathbf{x}^T \begin{pmatrix} 0 & \frac{1}{2} \\ \frac{1}{2} & 0 \end{pmatrix} \mathbf{x} = \mathbf{x}^T \mathbf{A} \mathbf{x}. \end{aligned}$$



Every quadratic form can be written in **standard** (i.e., diagonal) **form** since every real symmetric matrix is orthogonally similar to a diagonal matrix.

Example

Take

$$\begin{aligned} f(\mathbf{x}) &= x_1 x_2 = \mathbf{x}^T \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} \mathbf{x} \\ &= \mathbf{x}^T \begin{pmatrix} 0 & \frac{1}{2} \\ \frac{1}{2} & 0 \end{pmatrix} \mathbf{x} = \mathbf{x}^T \mathbf{A} \mathbf{x}. \end{aligned}$$

We want to find the standard form $f(\mathbf{y}) = \mathbf{y}^T \mathbf{D} \mathbf{y}$, where **D is diagonal** and **y are transformed coordinates**.



Example (cont.)

We can compute the eigenvalues and (orthogonal) eigenvectors of A , i.e.,

$$A = QDQ^T$$
$$\iff \frac{1}{2} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & -\frac{1}{2} \end{pmatrix} \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix}$$

Example (cont.)

We can compute the eigenvalues and (orthogonal) eigenvectors of A , i.e.,

$$A = QDQ^T$$

$$\iff \frac{1}{2} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & -\frac{1}{2} \end{pmatrix} \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix}$$

so that

$$f(\mathbf{x}) = \mathbf{x}^T A \mathbf{x} = \underbrace{\mathbf{x}^T Q}_{=\mathbf{y}^T} D Q^T \mathbf{x} = \mathbf{y}^T D \mathbf{y}$$

and the standard form is



Example (cont.)

We can compute the eigenvalues and (orthogonal) eigenvectors of A , i.e.,

$$A = QDQ^T$$

$$\iff \frac{1}{2} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & -\frac{1}{2} \end{pmatrix} \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix}$$

so that

$$f(\mathbf{x}) = \mathbf{x}^T A \mathbf{x} = \underbrace{\mathbf{x}^T Q}_{=\mathbf{y}^T} D Q^T \mathbf{x} = \mathbf{y}^T D \mathbf{y}$$

and the standard form is

$$\mathbf{y}^T \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & -\frac{1}{2} \end{pmatrix} \mathbf{y} = \frac{1}{2} (y_1^2 - y_2^2).$$



Remark

Instead of computing the eigenvalues and eigenvectors of A in the example, we can also consider the factorization

$$A = LDL^T.$$



Remark

Instead of computing the eigenvalues and eigenvectors of A in the example, we can also consider the factorization

$$A = LDL^T.$$

*For a positive definite A this is the **Cholesky factorization**, and it is cheaper to compute than eigenvalues and eigenvectors.*



Remark

Instead of computing the eigenvalues and eigenvectors of A in the example, we can also consider the factorization

$$A = LDL^T.$$

*For a positive definite A this is the **Cholesky factorization**, and it is cheaper to compute than eigenvalues and eigenvectors.*

Then

$$\mathbf{x}^T A \mathbf{x} = \underbrace{\mathbf{x}^T L}_{=\mathbf{y}^T} D L^T \mathbf{x} = \mathbf{y}^T D \mathbf{y} = \sum_{i=1}^n p_i y_i^2,$$

*where $D = \text{diag}(p_1, \dots, p_n)$ contains the **pivots used in Gaussian elimination**.*



Congruence transformations

Formally, the preceding argument uses a **congruence transformation**.

Definition

Two matrices $A, B \in \mathbb{R}^{n \times n}$ are called **congruent** if

$$B = C^T A C$$

for some nonsingular matrix C . Commonly used notation: $A \simeq B$.



Recall: A and B are **similar** if $B = P^{-1}AP$, and similar matrices have the same eigenvalues.



Recall: A and B are **similar** if $B = P^{-1}AP$, and similar matrices have the same eigenvalues.

Now,

Definition

Let A be a real symmetric matrix. The triple (ρ, ν, ζ) , where ρ , ν , and ζ , respectively, denote the number of **positive, negative, and zero eigenvalues of A** is called the **inertia** of A .



Theorem (Sylvester's Law of Inertia)

Let $A, B \in \mathbb{R}^{n \times n}$ be symmetric. Then A and B are *congruent*, i.e., $A \simeq B$, if and only if A and B have the *same inertias*.



Theorem (Sylvester's Law of Inertia)

Let $A, B \in \mathbb{R}^{n \times n}$ be symmetric. Then A and B are *congruent*, i.e., $A \simeq B$, if and only if A and B have the *same inertias*.

Proof.

See [Mey00].



Outline

- 1 Elementary Properties
- 2 Diagonalization via Similarity Transforms
- 3 Functions of Diagonalizable Matrices
- 4 Normal Matrices
- 5 Positive Definite Matrices
- 6 Iterative Solvers**
- 7 Krylov Methods



Iterative Solvers

Consider the linear system

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

where A has many zero entries, i.e., A is sparse.



Iterative Solvers

Consider the linear system

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

where A has many zero entries, i.e., A is sparse.

In this case, direct factorization methods (such as LU, QR, SVD) are very inefficient to solve $A\mathbf{x} = \mathbf{b}$.



Iterative Solvers

Consider the linear system

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

where A has many zero entries, i.e., A is sparse.

In this case, direct factorization methods (such as LU, QR, SVD) are very inefficient to solve $A\mathbf{x} = \mathbf{b}$.

Instead, one uses iterative solvers.



The **general framework for classical iterative solvers** is as follows:



The general framework for classical iterative solvers is as follows:

- We split A into

$$A = M - N,$$

where M^{-1} exists and — ideally — is easy to compute.



The **general framework for classical iterative solvers** is as follows:

- We **split A** into

$$A = M - N,$$

where M^{-1} **exists** and — ideally — **is easy to compute**.

- Then

$$Ax = b \iff (M - N)x = b$$



The **general framework for classical iterative solvers** is as follows:

- We **split A** into

$$A = M - N,$$

where M^{-1} **exists** and — ideally — **is easy to compute**.

- Then

$$Ax = b \iff (M - N)x = b \iff Mx = Nx + b$$



The **general framework for classical iterative solvers** is as follows:

- We **split** A into

$$A = M - N,$$

where M^{-1} **exists** and — ideally — **is easy to compute**.

- Then

$$A\mathbf{x} = \mathbf{b} \iff (M - N)\mathbf{x} = \mathbf{b} \iff M\mathbf{x} = N\mathbf{x} + \mathbf{b}$$

- and we **iterate**

$$M\mathbf{x}^{(k)} = N\mathbf{x}^{(k-1)} + \mathbf{b}$$



The **general framework for classical iterative solvers** is as follows:

- We **split A** into

$$A = M - N,$$

where M^{-1} **exists** and — ideally — **is easy to compute**.

- Then

$$A\mathbf{x} = \mathbf{b} \iff (M - N)\mathbf{x} = \mathbf{b} \iff M\mathbf{x} = N\mathbf{x} + \mathbf{b}$$

- and we **iterate**

$$\begin{aligned} M\mathbf{x}^{(k)} &= N\mathbf{x}^{(k-1)} + \mathbf{b} \\ \iff \mathbf{x}^{(k)} &= \underbrace{M^{-1}N}_{=H}\mathbf{x}^{(k-1)} + \underbrace{M^{-1}\mathbf{b}}_{=\mathbf{d}}, \quad k = 1, 2, 3, \dots, \end{aligned}$$

where $\mathbf{x}^{(0)}$ is some **initial guess** and $H = M^{-1}N$ is called the **iteration matrix**.



Theorem

Let M and N be two matrices such that $A = M - N$ and $H = M^{-1}N$. If $\rho(H) < 1$ then A is nonsingular and $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x} = A^{-1}\mathbf{b}$, i.e., the iterative method with iteration matrix H , converges for any initial guess $\mathbf{x}^{(0)}$ to the solution of $A\mathbf{x} = \mathbf{b}$.



Proof

First we show that A is nonsingular.



Proof

First we show that A is nonsingular.

Since $H = M^{-1}N$ (invertibility of M is an assumption) we have

$$\begin{aligned} A &= M - N \\ &= M - MH \\ &= M(I - H). \end{aligned} \tag{3}$$



Proof

First we show that A is nonsingular.

Since $H = M^{-1}N$ (invertibility of M is an assumption) we have

$$\begin{aligned} A &= M - N \\ &= M - MH \\ &= M(I - H). \end{aligned} \tag{3}$$

Now, since $\rho(H) < 1$ we know that $I - H$ is invertible via its Neumann series, and therefore A is invertible.



Proof (cont.)

Now we show that $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$:

Proof (cont.)

Now we show that $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$:

$$\mathbf{x}^{(k)} = \mathbf{H}\mathbf{x}^{(k-1)} + \mathbf{d}$$

Proof (cont.)

Now we show that $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$:

$$\begin{aligned} \mathbf{x}^{(k)} &= \mathbf{H}\mathbf{x}^{(k-1)} + \mathbf{d} \\ &= \mathbf{H} \left(\mathbf{H}\mathbf{x}^{(k-2)} + \mathbf{d} \right) + \mathbf{d} = \mathbf{H}^2 \mathbf{x}^{(k-2)} + (\mathbf{I} + \mathbf{H})\mathbf{d} \end{aligned}$$

Proof (cont.)

Now we show that $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$:

$$\begin{aligned}\mathbf{x}^{(k)} &= \mathbf{H}\mathbf{x}^{(k-1)} + \mathbf{d} \\ &= \mathbf{H} \left(\mathbf{H}\mathbf{x}^{(k-2)} + \mathbf{d} \right) + \mathbf{d} = \mathbf{H}^2\mathbf{x}^{(k-2)} + (\mathbf{I} + \mathbf{H})\mathbf{d} \\ &\vdots \\ &= \mathbf{H}^k\mathbf{x}^{(0)} + (\mathbf{I} + \mathbf{H} + \dots + \mathbf{H}^{k-1})\mathbf{d},\end{aligned}$$

where

Proof (cont.)

Now we show that $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$:

$$\begin{aligned}\mathbf{x}^{(k)} &= \mathbf{H}\mathbf{x}^{(k-1)} + \mathbf{d} \\ &= \mathbf{H} \left(\mathbf{H}\mathbf{x}^{(k-2)} + \mathbf{d} \right) + \mathbf{d} = \mathbf{H}^2\mathbf{x}^{(k-2)} + (\mathbf{I} + \mathbf{H})\mathbf{d} \\ &\vdots \\ &= \mathbf{H}^k\mathbf{x}^{(0)} + (\mathbf{I} + \mathbf{H} + \dots + \mathbf{H}^{k-1})\mathbf{d},\end{aligned}$$

where

$$\mathbf{H}^k \rightarrow \mathbf{O}$$

Proof (cont.)

Now we show that $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$:

$$\begin{aligned} \mathbf{x}^{(k)} &= \mathbf{H}\mathbf{x}^{(k-1)} + \mathbf{d} \\ &= \mathbf{H} \left(\mathbf{H}\mathbf{x}^{(k-2)} + \mathbf{d} \right) + \mathbf{d} = \mathbf{H}^2 \mathbf{x}^{(k-2)} + (\mathbf{I} + \mathbf{H})\mathbf{d} \\ &\vdots \\ &= \mathbf{H}^k \mathbf{x}^{(0)} + (\mathbf{I} + \mathbf{H} + \dots + \mathbf{H}^{k-1})\mathbf{d}, \end{aligned}$$

where

$$\mathbf{H}^k \rightarrow \mathbf{O} \quad \text{and} \quad (\mathbf{I} + \mathbf{H} + \dots + \mathbf{H}^{k-1}) \rightarrow (\mathbf{I} - \mathbf{H})^{-1} \quad \text{for } k \rightarrow \infty$$

Proof (cont.)

Now we show that $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$:

$$\begin{aligned} \mathbf{x}^{(k)} &= \mathbf{H}\mathbf{x}^{(k-1)} + \mathbf{d} \\ &= \mathbf{H} \left(\mathbf{H}\mathbf{x}^{(k-2)} + \mathbf{d} \right) + \mathbf{d} = \mathbf{H}^2 \mathbf{x}^{(k-2)} + (\mathbf{I} + \mathbf{H})\mathbf{d} \\ &\vdots \\ &= \mathbf{H}^k \mathbf{x}^{(0)} + (\mathbf{I} + \mathbf{H} + \dots + \mathbf{H}^{k-1})\mathbf{d}, \end{aligned}$$

where

$$\mathbf{H}^k \rightarrow \mathbf{O} \quad \text{and} \quad (\mathbf{I} + \mathbf{H} + \dots + \mathbf{H}^{k-1}) \rightarrow (\mathbf{I} - \mathbf{H})^{-1} \quad \text{for } k \rightarrow \infty$$

so that — using (3), i.e., $(\mathbf{I} - \mathbf{H})^{-1} = \mathbf{A}^{-1} \mathbf{M}$,

$$\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = (\mathbf{I} - \mathbf{H})^{-1} \mathbf{d}$$

Proof (cont.)

Now we show that $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$:

$$\begin{aligned} \mathbf{x}^{(k)} &= \mathbf{H}\mathbf{x}^{(k-1)} + \mathbf{d} \\ &= \mathbf{H} \left(\mathbf{H}\mathbf{x}^{(k-2)} + \mathbf{d} \right) + \mathbf{d} = \mathbf{H}^2 \mathbf{x}^{(k-2)} + (\mathbf{I} + \mathbf{H})\mathbf{d} \\ &\vdots \\ &= \mathbf{H}^k \mathbf{x}^{(0)} + (\mathbf{I} + \mathbf{H} + \dots + \mathbf{H}^{k-1})\mathbf{d}, \end{aligned}$$

where

$$\mathbf{H}^k \rightarrow \mathbf{O} \quad \text{and} \quad (\mathbf{I} + \mathbf{H} + \dots + \mathbf{H}^{k-1}) \rightarrow (\mathbf{I} - \mathbf{H})^{-1} \quad \text{for } k \rightarrow \infty$$

so that — using (3), i.e., $(\mathbf{I} - \mathbf{H})^{-1} = \mathbf{A}^{-1} \mathbf{M}$,

$$\begin{aligned} \lim_{k \rightarrow \infty} \mathbf{x}^{(k)} &= (\mathbf{I} - \mathbf{H})^{-1} \mathbf{d} \\ &= (\mathbf{I} - \mathbf{H})^{-1} \mathbf{M}^{-1} \mathbf{b} \end{aligned}$$

Proof (cont.)

Now we show that $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$:

$$\begin{aligned} \mathbf{x}^{(k)} &= \mathbf{H}\mathbf{x}^{(k-1)} + \mathbf{d} \\ &= \mathbf{H} \left(\mathbf{H}\mathbf{x}^{(k-2)} + \mathbf{d} \right) + \mathbf{d} = \mathbf{H}^2 \mathbf{x}^{(k-2)} + (\mathbf{I} + \mathbf{H})\mathbf{d} \\ &\vdots \\ &= \mathbf{H}^k \mathbf{x}^{(0)} + (\mathbf{I} + \mathbf{H} + \dots + \mathbf{H}^{k-1})\mathbf{d}, \end{aligned}$$

where

$$\mathbf{H}^k \rightarrow \mathbf{O} \quad \text{and} \quad (\mathbf{I} + \mathbf{H} + \dots + \mathbf{H}^{k-1}) \rightarrow (\mathbf{I} - \mathbf{H})^{-1} \quad \text{for } k \rightarrow \infty$$

so that — using (3), i.e., $(\mathbf{I} - \mathbf{H})^{-1} = \mathbf{A}^{-1}\mathbf{M}$,

$$\begin{aligned} \lim_{k \rightarrow \infty} \mathbf{x}^{(k)} &= (\mathbf{I} - \mathbf{H})^{-1} \mathbf{d} \\ &= (\mathbf{I} - \mathbf{H})^{-1} \mathbf{M}^{-1} \mathbf{b} = \mathbf{A}^{-1} \mathbf{b} \end{aligned}$$

Proof (cont.)

Now we show that $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$:

$$\begin{aligned} \mathbf{x}^{(k)} &= \mathbf{H} \mathbf{x}^{(k-1)} + \mathbf{d} \\ &= \mathbf{H} \left(\mathbf{H} \mathbf{x}^{(k-2)} + \mathbf{d} \right) + \mathbf{d} = \mathbf{H}^2 \mathbf{x}^{(k-2)} + (\mathbf{I} + \mathbf{H}) \mathbf{d} \\ &\vdots \\ &= \mathbf{H}^k \mathbf{x}^{(0)} + (\mathbf{I} + \mathbf{H} + \dots + \mathbf{H}^{k-1}) \mathbf{d}, \end{aligned}$$

where

$$\mathbf{H}^k \rightarrow \mathbf{O} \quad \text{and} \quad (\mathbf{I} + \mathbf{H} + \dots + \mathbf{H}^{k-1}) \rightarrow (\mathbf{I} - \mathbf{H})^{-1} \quad \text{for } k \rightarrow \infty$$

so that — using (3), i.e., $(\mathbf{I} - \mathbf{H})^{-1} = \mathbf{A}^{-1} \mathbf{M}$,

$$\begin{aligned} \lim_{k \rightarrow \infty} \mathbf{x}^{(k)} &= (\mathbf{I} - \mathbf{H})^{-1} \mathbf{d} \\ &= (\mathbf{I} - \mathbf{H})^{-1} \mathbf{M}^{-1} \mathbf{b} = \mathbf{A}^{-1} \mathbf{b} = \mathbf{x}. \end{aligned}$$

Remark

In order to have a “good” iterative solver we will want



Remark

In order to have a “good” iterative solver we will want

- *fast convergence* — ensured by $\rho(H) \ll 1$,



Remark

In order to have a “good” iterative solver we will want

- *fast convergence* — ensured by $\rho(H) \ll 1$,
- *simple computation* — ensured by easy computation of M^{-1} (or $H = M^{-1}N$ and $\mathbf{d} = M^{-1}\mathbf{b}$).



Remark

In order to have a “good” iterative solver we will want

- *fast convergence* — ensured by $\rho(H) \ll 1$,
- *simple computation* — ensured by easy computation of M^{-1} (or $H = M^{-1}N$ and $\mathbf{d} = M^{-1}\mathbf{b}$).

We conclude by presenting two standard examples:

- *Jacobi iteration,*
- *Gauss-Seidel iteration.*



Jacobi iteration

We take $M = D = \text{diag}(A)$, which is easy to invert.



Jacobi iteration

We take $M = D = \text{diag}(A)$, which is easy to invert.

Then

$$A = M - N = D - N,$$

i.e., $N = -(A - D)$ or, if $A = L + D + U$, $N = -(L + U)$.



Jacobi iteration

We take $M = D = \text{diag}(A)$, which is easy to invert.

Then

$$A = M - N = D - N,$$

i.e., $N = -(A - D)$ or, if $A = L + D + U$, $N = -(L + U)$.

Therefore $A\mathbf{x} = \mathbf{b}$ is solved via

$$D\mathbf{x}^{(k)} = N\mathbf{x}^{(k-1)} + \mathbf{b}, \quad k = 1, 2, 3, \dots,$$



Jacobi iteration

We take $M = D = \text{diag}(A)$, which is easy to invert.

Then

$$A = M - N = D - N,$$

i.e., $N = -(A - D)$ or, if $A = L + D + U$, $N = -(L + U)$.

Therefore $A\mathbf{x} = \mathbf{b}$ is solved via

$$D\mathbf{x}^{(k)} = N\mathbf{x}^{(k-1)} + \mathbf{b}, \quad k = 1, 2, 3, \dots,$$

or componentwise

$$\mathbf{x}_i^{(k)} = \frac{1}{a_{ii}} \left(b_i - \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} \mathbf{x}_j^{(k-1)} \right), \quad i = 1, 2, \dots, n.$$



Remark

- *Jacobi iteration is **embarrassingly parallel**, i.e., the above componentwise equations can be directly implemented on n parallel processors.*



Remark

- *Jacobi iteration is **embarrassingly parallel**, i.e., the above componentwise equations can be directly implemented on n parallel processors.*
- *Also, **only entries from the i^{th} row** of the matrix are **needed to update the i^{th} component of \mathbf{x} .***



Remark

- *Jacobi iteration is **embarrassingly parallel**, i.e., the above componentwise equations can be directly implemented on n parallel processors.*
- *Also, **only entries from the i^{th} row** of the matrix are **needed to update the i^{th} component of \mathbf{x} .***
- *Jacobi iteration had long been considered as **too simple (and too slow) to be useful**. However, a **recent modification** [YM14] using relaxation **has changed that**. This modification was customized to solve elliptic PDEs via a finite difference discretization.*



Theorem

If A is diagonally dominant, then Jacobi iteration converges for any initial guess.



Proof.

Diagonal dominance says

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad i = 1, \dots, n$$

Proof.

Diagonal dominance says

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad i = 1, \dots, n \iff \sum_{\substack{j=1 \\ j \neq i}}^n \left| \frac{a_{ij}}{a_{ii}} \right| < 1.$$

Proof.

Diagonal dominance says

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad i = 1, \dots, n \iff \sum_{\substack{j=1 \\ j \neq i}}^n \left| \frac{a_{ij}}{a_{ii}} \right| < 1.$$

Now

$$\|H\|_{\infty} =$$



Proof.

Diagonal dominance says

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad i = 1, \dots, n \iff \sum_{\substack{j=1 \\ j \neq i}}^n \left| \frac{a_{ij}}{a_{ii}} \right| < 1.$$

Now

$$\|H\|_{\infty} = \|D^{-1}N\|_{\infty}$$



Proof.

Diagonal dominance says

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad i = 1, \dots, n \iff \sum_{\substack{j=1 \\ j \neq i}}^n \left| \frac{a_{ij}}{a_{ii}} \right| < 1.$$

Now

$$\|H\|_{\infty} = \|D^{-1}N\|_{\infty} = \max_{1 \leq i \leq n} \sum_{j=1}^n \left| \frac{a_{ij}}{a_{ii}} \right|$$



Proof.

Diagonal dominance says

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad i = 1, \dots, n \iff \sum_{\substack{j=1 \\ j \neq i}}^n \left| \frac{a_{ij}}{a_{ii}} \right| < 1.$$

Now

$$\|H\|_{\infty} = \|D^{-1}N\|_{\infty} = \max_{1 \leq i \leq n} \sum_{j=1}^n \left| \frac{a_{ij}}{a_{ii}} \right|$$

$$\stackrel{\text{diag}(N)=0}{=} \max_{1 \leq i \leq n} \sum_{\substack{j=1 \\ j \neq i}}^n \left| \frac{a_{ij}}{a_{ii}} \right| < 1.$$



Remark

Since $\rho(H) < \|H\|$, *diagonal dominance (or $\|H\|_\infty < 1$) is a weaker condition than $\rho(H) < 1$.*



Gauss–Seidel iteration

Let's again decompose $A = L + D + U$, but now **take**

$$M = D + L, \quad N = -U.$$



Gauss–Seidel iteration

Let's again decompose $A = L + D + U$, but now **take**

$$M = D + L, \quad N = -U.$$

Then

$$H = M^{-1}N = -(D + L)^{-1}U$$

$$\mathbf{d} = M^{-1}\mathbf{b} = (D + L)^{-1}\mathbf{b}.$$



Gauss–Seidel iteration

Let's again decompose $A = L + D + U$, but now **take**

$$M = D + L, \quad N = -U.$$

Then

$$H = M^{-1}N = -(D + L)^{-1}U$$

$$d = M^{-1}b = (D + L)^{-1}b.$$

The **iteration formula** is

$$x^{(k)} = -(D + L)^{-1}Ux^{(k-1)} + (D + L)^{-1}b$$



Gauss–Seidel iteration

Let's again decompose $A = L + D + U$, but now **take**

$$M = D + L, \quad N = -U.$$

Then

$$H = M^{-1}N = -(D + L)^{-1}U$$

$$\mathbf{d} = M^{-1}\mathbf{b} = (D + L)^{-1}\mathbf{b}.$$

The **iteration formula** is

$$\begin{aligned} \mathbf{x}^{(k)} &= -(D + L)^{-1}U\mathbf{x}^{(k-1)} + (D + L)^{-1}\mathbf{b} \\ \iff (D + L)\mathbf{x}^{(k)} &= \mathbf{b} - U\mathbf{x}^{(k-1)}. \end{aligned}$$



Componentwise we get

$$\sum_{j=1}^{i-1} a_{ij} \mathbf{x}_j^{(k)} + a_{ii} \mathbf{x}_i^{(k)} = b_i - \sum_{j=i+1}^n a_{ij} \mathbf{x}_j^{(k-1)}.$$



Componentwise we get

$$\sum_{j=1}^{i-1} a_{ij} \mathbf{x}_j^{(k)} + a_{ii} \mathbf{x}_i^{(k)} = b_i - \sum_{j=i+1}^n a_{ij} \mathbf{x}_j^{(k-1)}.$$

Since when we work on the i^{th} component the components $\mathbf{x}_j^{(k)}$, $j < i$, have already been updated we can write



Componentwise we get

$$\sum_{j=1}^{i-1} a_{ij} \mathbf{x}_j^{(k)} + a_{ii} \mathbf{x}_i^{(k)} = b_i - \sum_{j=i+1}^n a_{ij} \mathbf{x}_j^{(k-1)}.$$

Since when we work on the i^{th} component the components $\mathbf{x}_j^{(k)}$, $j < i$, have already been updated we can write

$$\mathbf{x}_i^{(k)} = \frac{1}{a_{ii}} \left(b_i - \sum_{j=1}^{i-1} a_{ij} \mathbf{x}_j^{(k)} - \sum_{j=i+1}^n a_{ij} \mathbf{x}_j^{(k-1)} \right), \quad i = 1, 2, \dots, n.$$



Componentwise we get

$$\sum_{j=1}^{i-1} a_{ij} \mathbf{x}_j^{(k)} + a_{ii} \mathbf{x}_i^{(k)} = b_i - \sum_{j=i+1}^n a_{ij} \mathbf{x}_j^{(k-1)}.$$

Since when we work on the i^{th} component the components $\mathbf{x}_j^{(k)}$, $j < i$, have already been updated we can write

$$\mathbf{x}_i^{(k)} = \frac{1}{a_{ii}} \left(b_i - \sum_{j=1}^{i-1} a_{ij} \mathbf{x}_j^{(k)} - \sum_{j=i+1}^n a_{ij} \mathbf{x}_j^{(k-1)} \right), \quad i = 1, 2, \dots, n.$$

Remark

Gauss–Seidel iteration is similar to Jacobi iteration, but it uses the most recently computed information as soon as it becomes available (instead of waiting until the next iteration, as Jacobi does).

Convergence of Gauss-Seidel iteration

Theorem

Gauss–Seidel iteration converges for any initial guess if

- 1 A is *diagonally dominant*, or
- 2 A is *symmetric positive definite*.



Convergence of Gauss-Seidel iteration

Theorem

Gauss–Seidel iteration converges for any initial guess if

- 1 *A is diagonally dominant, or*
- 2 *A is symmetric positive definite.*

Proof.

- 1 In [Mey00],
- 2 on next few slides.



Convergence of Gauss-Seidel iteration

Theorem

Gauss–Seidel iteration converges for any initial guess if

- 1 A is *diagonally dominant*, or
- 2 A is *symmetric positive definite*.

Proof.

- 1 In [Mey00],
- 2 on next few slides.



Remark

Usually Gauss–Seidel converges faster than Jacobi. However, there are exceptions.

Proof (convergence for positive definite A)

Since A is symmetric, we can decompose

$$A = L + D + L^T, \quad H = -(D + L)^{-1}L^T.$$

Proof (convergence for positive definite A)

Since **A is symmetric**, we can **decompose**

$$A = L + D + L^T, \quad H = -(D + L)^{-1}L^T.$$

Convergence will be ensured if we can **show that $\rho(H) < 1$** , i.e.,

$$\rho(-(D + L)^{-1}L^T) < 1.$$

Proof (convergence for positive definite A)

Since **A is symmetric**, we can **decompose**

$$A = L + D + L^T, \quad H = -(D + L)^{-1}L^T.$$

Convergence will be ensured if we can **show that $\rho(H) < 1$** , i.e.,

$$\rho(-(D + L)^{-1}L^T) < 1.$$

Since **D** has positive entries (otherwise A couldn't be positive definite), D is positive definite (and therefore **nonsingular**) so that

$$\tilde{H} = D^{1/2}HD^{-1/2}$$

has the **same eigenvalues** as H.

Proof (convergence for positive definite A)

Since **A is symmetric**, we can **decompose**

$$A = L + D + L^T, \quad H = -(D + L)^{-1}L^T.$$

Convergence will be ensured if we can **show that $\rho(H) < 1$** , i.e.,

$$\rho(-(D + L)^{-1}L^T) < 1.$$

Since **D** has positive entries (otherwise A couldn't be positive definite), D is positive definite (and therefore **nonsingular**) so that

$$\tilde{H} = D^{1/2}HD^{-1/2}$$

has the **same eigenvalues** as H.

Therefore, we now **show that**

$$\rho(\tilde{H}) < 1.$$

Proof (cont.)

First, we rewrite \tilde{H} . For this we require a **push-through identity for the matrix inverse** ([Ber09], similar to what we had in Chapter 3):

$$(I + AB)^{-1}A = A(I + BA)^{-1}. \quad (4)$$

Proof (cont.)

First, we rewrite \tilde{H} . For this we require a **push-through identity for the matrix inverse** ([Ber09], similar to what we had in Chapter 3):

$$(I + AB)^{-1}A = A(I + BA)^{-1}. \quad (4)$$

If we let $A = D^{-1/2}$ and $B = LD^{-1/2}$, then we get

$$(I + D^{-1/2}LD^{-1/2})^{-1}D^{-1/2} \stackrel{(4)}{=} D^{-1/2}(I + LD^{-1/2}D^{-1/2})^{-1}$$

(5)



Proof (cont.)

First, we rewrite \tilde{H} . For this we require a **push-through identity for the matrix inverse** ([Ber09], similar to what we had in Chapter 3):

$$(I + AB)^{-1}A = A(I + BA)^{-1}. \quad (4)$$

If we let $A = D^{-1/2}$ and $B = LD^{-1/2}$, then we get

$$\begin{aligned} (I + D^{-1/2}LD^{-1/2})^{-1}D^{-1/2} &\stackrel{(4)}{=} D^{-1/2}(I + LD^{-1/2}D^{-1/2})^{-1} \\ &= D^{-1/2}(I + LD^{-1})^{-1} \end{aligned}$$

(5)

Proof (cont.)

First, we rewrite \tilde{H} . For this we require a **push-through identity for the matrix inverse** ([Ber09], similar to what we had in Chapter 3):

$$(I + AB)^{-1}A = A(I + BA)^{-1}. \quad (4)$$

If we let $A = D^{-1/2}$ and $B = LD^{-1/2}$, then we get

$$\begin{aligned} (I + D^{-1/2}LD^{-1/2})^{-1}D^{-1/2} &\stackrel{(4)}{=} D^{-1/2}(I + LD^{-1/2}D^{-1/2})^{-1} \\ &= D^{-1/2}(I + LD^{-1})^{-1} \\ &= D^{-1/2}(DD^{-1} + LD^{-1})^{-1} \end{aligned}$$

(5)



Proof (cont.)

First, we rewrite \tilde{H} . For this we require a **push-through identity for the matrix inverse** ([Ber09], similar to what we had in Chapter 3):

$$(I + AB)^{-1}A = A(I + BA)^{-1}. \quad (4)$$

If we let $A = D^{-1/2}$ and $B = LD^{-1/2}$, then we get

$$\begin{aligned} (I + D^{-1/2}LD^{-1/2})^{-1}D^{-1/2} &\stackrel{(4)}{=} D^{-1/2}(I + LD^{-1/2}D^{-1/2})^{-1} \\ &= D^{-1/2}(I + LD^{-1})^{-1} \\ &= D^{-1/2}(DD^{-1} + LD^{-1})^{-1} \\ &= D^{-1/2} \left((D + L)D^{-1} \right)^{-1} \end{aligned} \quad (5)$$

Proof (cont.)

First, we rewrite \tilde{H} . For this we require a **push-through identity for the matrix inverse** ([Ber09], similar to what we had in Chapter 3):

$$(I + AB)^{-1}A = A(I + BA)^{-1}. \quad (4)$$

If we let $A = D^{-1/2}$ and $B = LD^{-1/2}$, then we get

$$\begin{aligned} (I + D^{-1/2}LD^{-1/2})^{-1}D^{-1/2} &\stackrel{(4)}{=} D^{-1/2}(I + LD^{-1/2}D^{-1/2})^{-1} \\ &= D^{-1/2}(I + LD^{-1})^{-1} \\ &= D^{-1/2}(DD^{-1} + LD^{-1})^{-1} \\ &= D^{-1/2}((D + L)D^{-1})^{-1} \\ &= D^{-1/2}D(D + L)^{-1} = D^{1/2}(D + L)^{-1}. \quad (5) \end{aligned}$$



Proof (cont.)

Therefore

$$\begin{aligned}\tilde{H} &= D^{1/2}HD^{-1/2} \\ &= -D^{1/2}(D+L)^{-1}L^T D^{-1/2}\end{aligned}$$

Proof (cont.)

Therefore

$$\begin{aligned}\tilde{H} &= D^{1/2}HD^{-1/2} \\ &= -D^{1/2}(D+L)^{-1}L^T D^{-1/2} \\ &\stackrel{(5)}{=} -(I + D^{-1/2}LD^{-1/2})^{-1}D^{-1/2}L^T D^{-1/2}\end{aligned}$$

Proof (cont.)

Therefore

$$\begin{aligned}
 \tilde{H} &= D^{1/2}HD^{-1/2} \\
 &= -D^{1/2}(D+L)^{-1}L^T D^{-1/2} \\
 &\stackrel{(5)}{=} -(I + D^{-1/2}LD^{-1/2})^{-1}D^{-1/2}L^T D^{-1/2} \\
 &= -(I + \tilde{L})^{-1}\tilde{L}^T,
 \end{aligned}$$

where $\tilde{L} = D^{-1/2}LD^{-1/2}$.

Proof (cont.)

Therefore

$$\begin{aligned}
 \tilde{H} &= D^{1/2}HD^{-1/2} \\
 &= -D^{1/2}(D + L)^{-1}L^T D^{-1/2} \\
 &\stackrel{(5)}{=} -(I + D^{-1/2}LD^{-1/2})^{-1}D^{-1/2}L^T D^{-1/2} \\
 &= -(I + \tilde{L})^{-1}\tilde{L}^T,
 \end{aligned}$$

where $\tilde{L} = D^{-1/2}LD^{-1/2}$.

Now consider an **eigenpair** (λ, \mathbf{x}) of \tilde{H} with $\mathbf{x}^* \mathbf{x} = 1$.

Proof (cont.)

Therefore

$$\begin{aligned}
 \tilde{H} &= D^{1/2}HD^{-1/2} \\
 &= -D^{1/2}(D+L)^{-1}L^TD^{-1/2} \\
 &\stackrel{(5)}{=} -(I+D^{-1/2}LD^{-1/2})^{-1}D^{-1/2}L^TD^{-1/2} \\
 &= -(I+\tilde{L})^{-1}\tilde{L}^T,
 \end{aligned}$$

where $\tilde{L} = D^{-1/2}LD^{-1/2}$.

Now consider an **eigenpair** (λ, \mathbf{x}) of \tilde{H} with $\mathbf{x}^*\mathbf{x} = 1$. Then

$$\tilde{H}\mathbf{x} = \lambda\mathbf{x} \iff -\tilde{L}^T\mathbf{x} = \lambda(I+\tilde{L})\mathbf{x}.$$

Proof (cont.)

Therefore

$$\begin{aligned}
 \tilde{H} &= D^{1/2} H D^{-1/2} \\
 &= -D^{1/2} (D + L)^{-1} L^T D^{-1/2} \\
 &\stackrel{(5)}{=} -(I + D^{-1/2} L D^{-1/2})^{-1} D^{-1/2} L^T D^{-1/2} \\
 &= -(I + \tilde{L})^{-1} \tilde{L}^T,
 \end{aligned}$$

where $\tilde{L} = D^{-1/2} L D^{-1/2}$.Now consider an eigenpair (λ, \mathbf{x}) of \tilde{H} with $\mathbf{x}^* \mathbf{x} = 1$. Then

$$\tilde{H} \mathbf{x} = \lambda \mathbf{x} \iff -\tilde{L}^T \mathbf{x} = \lambda (I + \tilde{L}) \mathbf{x}.$$

Multiplying by \mathbf{x}^* yields

$$-\mathbf{x}^* \tilde{L}^T \mathbf{x} = \lambda (\underbrace{\mathbf{x}^* \mathbf{x}}_{=1} + \mathbf{x}^* \tilde{L} \mathbf{x}) \iff \lambda = \frac{-\mathbf{x}^* \tilde{L}^T \mathbf{x}}{1 + \mathbf{x}^* \tilde{L} \mathbf{x}}.$$

Proof (cont.)

Finally, we let $\mathbf{x}^* \tilde{\mathbf{L}} \mathbf{x} = a + bi$.



Proof (cont.)

Finally, we let $\mathbf{x}^* \tilde{\mathbf{L}} \mathbf{x} = a + bi$. Then we have $\mathbf{x}^* \tilde{\mathbf{L}}^T \mathbf{x} = a - bi$ so that

$$|\lambda|^2 = \left| \frac{-a + bi}{1 + a + bi} \right|^2$$



Proof (cont.)

Finally, we let $\mathbf{x}^* \tilde{\mathbf{L}} \mathbf{x} = a + bi$. Then we have $\mathbf{x}^* \tilde{\mathbf{L}}^T \mathbf{x} = a - bi$ so that

$$|\lambda|^2 = \left| \frac{-a + bi}{1 + a + bi} \right|^2 = \frac{a^2 + b^2}{1 + 2a + a^2 + b^2}$$



Proof (cont.)

Finally, we let $\mathbf{x}^* \tilde{\mathbf{L}} \mathbf{x} = a + bi$. Then we have $\mathbf{x}^* \tilde{\mathbf{L}}^T \mathbf{x} = a - bi$ so that

$$|\lambda|^2 = \left| \frac{-a + bi}{1 + a + bi} \right|^2 = \frac{a^2 + b^2}{1 + 2a + a^2 + b^2} < 1$$

since $1 + 2a > 0$, as we now show:



Proof (cont.)

Finally, we let $\mathbf{x}^* \tilde{\mathbf{L}} \mathbf{x} = a + bi$. Then we have $\mathbf{x}^* \tilde{\mathbf{L}}^T \mathbf{x} = a - bi$ so that

$$|\lambda|^2 = \left| \frac{-a + bi}{1 + a + bi} \right|^2 = \frac{a^2 + b^2}{1 + 2a + a^2 + b^2} < 1$$

since $1 + 2a > 0$, as we now show:

The matrix $\mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2} = \tilde{\mathbf{L}} + \mathbf{I} + \tilde{\mathbf{L}}^T$ is **positive definite**, and therefore its **quadratic form is positive**.



Proof (cont.)

Finally, we let $\mathbf{x}^* \tilde{\mathbf{L}} \mathbf{x} = a + bi$. Then we have $\mathbf{x}^* \tilde{\mathbf{L}}^T \mathbf{x} = a - bi$ so that

$$|\lambda|^2 = \left| \frac{-a + bi}{1 + a + bi} \right|^2 = \frac{a^2 + b^2}{1 + 2a + a^2 + b^2} < 1$$

since $1 + 2a > 0$, as we now show:

The matrix $\mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2} = \tilde{\mathbf{L}} + \mathbf{I} + \tilde{\mathbf{L}}^T$ is **positive definite**, and therefore its **quadratic form is positive**.

In particular, **using the eigenvector \mathbf{x}** we have

$$0 < \mathbf{x}^* \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2} \mathbf{x}$$



Proof (cont.)

Finally, we let $\mathbf{x}^* \tilde{\mathbf{L}} \mathbf{x} = a + bi$. Then we have $\mathbf{x}^* \tilde{\mathbf{L}}^T \mathbf{x} = a - bi$ so that

$$|\lambda|^2 = \left| \frac{-a + bi}{1 + a + bi} \right|^2 = \frac{a^2 + b^2}{1 + 2a + a^2 + b^2} < 1$$

since $1 + 2a > 0$, as we now show:

The matrix $\mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2} = \tilde{\mathbf{L}} + \mathbf{I} + \tilde{\mathbf{L}}^T$ is **positive definite**, and therefore its **quadratic form is positive**.

In particular, **using the eigenvector \mathbf{x}** we have

$$0 < \mathbf{x}^* \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2} \mathbf{x} = \underbrace{\mathbf{x}^* \tilde{\mathbf{L}} \mathbf{x}}_{=a+bi} + \underbrace{\mathbf{x}^* \mathbf{x}}_{=1} + \underbrace{\mathbf{x}^* \tilde{\mathbf{L}}^T \mathbf{x}}_{=a-bi}$$



Proof (cont.)

Finally, we let $\mathbf{x}^* \tilde{\mathbf{L}} \mathbf{x} = a + bi$. Then we have $\mathbf{x}^* \tilde{\mathbf{L}}^T \mathbf{x} = a - bi$ so that

$$|\lambda|^2 = \left| \frac{-a + bi}{1 + a + bi} \right|^2 = \frac{a^2 + b^2}{1 + 2a + a^2 + b^2} < 1$$

since $1 + 2a > 0$, as we now show:

The matrix $\mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2} = \tilde{\mathbf{L}} + \mathbf{I} + \tilde{\mathbf{L}}^T$ is **positive definite**, and therefore its **quadratic form is positive**.

In particular, **using the eigenvector \mathbf{x}** we have

$$0 < \mathbf{x}^* \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2} \mathbf{x} = \underbrace{\mathbf{x}^* \tilde{\mathbf{L}} \mathbf{x}}_{=a+bi} + \underbrace{\mathbf{x}^* \mathbf{x}}_{=1} + \underbrace{\mathbf{x}^* \tilde{\mathbf{L}}^T \mathbf{x}}_{=a-bi} = 1 + 2a.$$



Outline

- 1 Elementary Properties
- 2 Diagonalization via Similarity Transforms
- 3 Functions of Diagonalizable Matrices
- 4 Normal Matrices
- 5 Positive Definite Matrices
- 6 Iterative Solvers
- 7 Krylov Methods**



Krylov Methods

We end with a very brief overview of Krylov methods.

This class of methods includes many of the **state-of-the-art numerical methods** for solving

$$A\mathbf{x} = \mathbf{b} \quad \text{or} \quad A\mathbf{x} = \lambda\mathbf{x}.$$

Some examples include:

- **Linear system solvers:**
 - conjugate gradient (CG), biconjugate gradient (BiCG), biconjugate gradient stabilized (BiCGSTAB), minimal residual (MINRES), generalized minimum residual (GMRES)
- **Eigensolvers:**
 - Lanczos iteration, Arnoldi iteration



The **basic building blocks** for all these methods are

Definition

For an $n \times n$ matrix A and nonzero n -vector \mathbf{b} we define

Krylov sequence: $\{\mathbf{b}, A\mathbf{b}, A^2\mathbf{b}, \dots\}$,

Krylov subspace: $\mathcal{K}_j = \text{span}\{\mathbf{b}, A\mathbf{b}, \dots, A^{j-1}\mathbf{b}\}$,

Krylov matrix: $K = (\mathbf{b} \quad A\mathbf{b} \quad \dots \quad A^{j-1}\mathbf{b})$.



Consider

$$\begin{aligned}AK &= (\mathbf{Ab} \quad \mathbf{A}^2\mathbf{b} \quad \cdots \quad \mathbf{A}^j\mathbf{b}) \\ &= \mathbf{K} (\mathbf{e}_2 \quad \mathbf{e}_3 \quad \cdots \quad \mathbf{e}_j \quad -\mathbf{c}),\end{aligned}$$

where $\mathbf{c} = -\mathbf{K}^{-1}\mathbf{A}^j\mathbf{b}$.

Note that the first $j - 1$ columns of AK coincide with columns 2 to j of \mathbf{K} .



Consider

$$\begin{aligned} AK &= (\mathbf{A}\mathbf{b} \quad \mathbf{A}^2\mathbf{b} \quad \cdots \quad \mathbf{A}^j\mathbf{b}) \\ &= K (\mathbf{e}_2 \quad \mathbf{e}_3 \quad \cdots \quad \mathbf{e}_j \quad -\mathbf{c}), \end{aligned}$$

where $\mathbf{c} = -K^{-1}\mathbf{A}^j\mathbf{b}$.

Note that the **first $j - 1$ columns of AK coincide with columns 2 to j of K** .

Letting $\mathbf{C} = (\mathbf{e}_2 \quad \mathbf{e}_3 \quad \cdots \quad \mathbf{e}_j \quad -\mathbf{c})$ we therefore have

$$AK = KC \iff K^{-1}AK = \mathbf{C},$$

i.e., \mathbf{A} and \mathbf{C} are **similar** and have the **same eigenvalues**.



Consider

$$\begin{aligned} AK &= (A\mathbf{b} \quad A^2\mathbf{b} \quad \dots \quad A^j\mathbf{b}) \\ &= K (\mathbf{e}_2 \quad \mathbf{e}_3 \quad \dots \quad \mathbf{e}_j \quad -\mathbf{c}), \end{aligned}$$

where $\mathbf{c} = -K^{-1}A^j\mathbf{b}$.

Note that the **first $j - 1$ columns of AK coincide with columns 2 to j of K** .

Letting $C = (\mathbf{e}_2 \quad \mathbf{e}_3 \quad \dots \quad \mathbf{e}_j \quad -\mathbf{c})$ we therefore have

$$AK = KC \iff K^{-1}AK = C,$$

i.e., A and C are **similar** and have the **same eigenvalues**.

Remark

The matrix C is called a **companion matrix**. It is **upper Hessenberg**, i.e., upper triangular with an additional nonzero subdiagonal.

Computation with such matrices can be performed quite efficiently.

If $j = n$ and we use **exact arithmetic** then $\mathcal{K}_n = R(\mathbf{A})$.



If $j = n$ and we use **exact arithmetic** then $\mathcal{K}_n = R(A)$.

Since we know that $\mathbf{x} \in R(A)$, the **fundamental idea of a Krylov method** is to

- iteratively produce approximate solutions \mathbf{x}_j that are projections into \mathcal{K}_j
- with the hope that low-dimensional Krylov subspaces already contain most of the essential information about $R(A)$.



If $j = n$ and we use **exact arithmetic** then $\mathcal{K}_n = R(A)$.

Since we know that $\mathbf{x} \in R(A)$, the **fundamental idea of a Krylov method** is to

- iteratively produce approximate solutions \mathbf{x}_j that are projections into \mathcal{K}_j
- with the hope that low-dimensional Krylov subspaces already contain most of the essential information about $R(A)$.

The main **practical problem with Krylov subspaces** is that the vectors $A^j \mathbf{b}$ all approach the dominant eigenvector of A (cf. **power method**), and so the **Krylov matrix K becomes ill-conditioned**.



The goal of all Krylov methods now is to find better bases for the Krylov subspaces \mathcal{K}_j .



The goal of all Krylov methods now is to find better bases for the Krylov subspaces \mathcal{K}_j .

This is essentially done via QR factorization, i.e., $K = QR$ leads to

$$\begin{aligned} AK = KC &\iff AQR = QRC \\ &\iff Q^T A Q = R C R^{-1} = H, \end{aligned}$$

where H is another upper Hessenberg matrix.



Arnoldi iteration

Arnoldi iteration is the **standard algorithm used to find the matrices Q and H .**

At the j^{th} iteration it will produce matrices

- Q_j , $n \times j$ with orthogonal columns that form a basis for \mathcal{K}_j ;
- Q_{j+1} , $n \times j + 1$ with orthogonal columns that form a basis for \mathcal{K}_{j+1} ;
- \tilde{H}_j , upper Hessenberg.

These matrices satisfy

$$AQ_j = Q_{j+1}\tilde{H}_j.$$



GMRES

The GMRES method attempts to solve $Ax = b$ by minimizing the residual $\|b - Ax_j\|_2$ at each iteration.



GMRES

The GMRES methods attempts to solve $A\mathbf{x} = \mathbf{b}$ by minimizing the residual $\|\mathbf{b} - A\mathbf{x}_j\|_2$ at each iteration.

Since the approximate solution $\mathbf{x}_j \in \mathcal{K}_j$ we can express it using an orthogonal basis, i.e.,

$$\mathbf{x}_j = Q_j \mathbf{z},$$

for an appropriate \mathbf{z} .



GMRES

The GMRES methods attempts to solve $A\mathbf{x} = \mathbf{b}$ by minimizing the residual $\|\mathbf{b} - A\mathbf{x}_j\|_2$ at each iteration.

Since the approximate solution $\mathbf{x}_j \in \mathcal{K}_j$ we can express it using an orthogonal basis, i.e.,

$$\mathbf{x}_j = Q_j \mathbf{z},$$

for an appropriate \mathbf{z} .

Then

$$\|\mathbf{b} - A\mathbf{x}_j\|_2 = \|\mathbf{b} - AQ_j\mathbf{z}\|_2 = \|\mathbf{b} - Q_{j+1}\tilde{H}_j\mathbf{z}\|_2.$$



GMRES

The GMRES methods attempts to solve $\mathbf{Ax} = \mathbf{b}$ by minimizing the residual $\|\mathbf{b} - \mathbf{Ax}_j\|_2$ at each iteration.

Since the approximate solution $\mathbf{x}_j \in \mathcal{K}_j$ we can express it using an orthogonal basis, i.e.,

$$\mathbf{x}_j = \mathbf{Q}_j \mathbf{z},$$

for an appropriate \mathbf{z} .

Then

$$\|\mathbf{b} - \mathbf{Ax}_j\|_2 = \|\mathbf{b} - \mathbf{AQ}_j \mathbf{z}\|_2 = \|\mathbf{b} - \mathbf{Q}_{j+1} \tilde{\mathbf{H}}_j \mathbf{z}\|_2.$$

Multiplication by an orthogonal matrix does not change the 2-norm, so

$$\|\mathbf{b} - \mathbf{Ax}_j\|_2 = \|\underbrace{\mathbf{Q}_{j+1}^T \mathbf{Q}_{j+1}}_{=I} \mathbf{b} - \tilde{\mathbf{H}}_j \mathbf{z}\|_2.$$



GMRES

The GMRES methods attempts to solve $A\mathbf{x} = \mathbf{b}$ by minimizing the residual $\|\mathbf{b} - A\mathbf{x}_j\|_2$ at each iteration.

Since the approximate solution $\mathbf{x}_j \in \mathcal{K}_j$ we can express it using an orthogonal basis, i.e.,

$$\mathbf{x}_j = Q_j \mathbf{z},$$

for an appropriate \mathbf{z} .

Then

$$\|\mathbf{b} - A\mathbf{x}_j\|_2 = \|\mathbf{b} - AQ_j\mathbf{z}\|_2 = \|\mathbf{b} - Q_{j+1}\tilde{H}_j\mathbf{z}\|_2.$$

Multiplication by an orthogonal matrix does not change the 2-norm, so

$$\|\mathbf{b} - A\mathbf{x}_j\|_2 = \|Q_{j+1}^T \mathbf{b} - \underbrace{Q_{j+1}^T Q_{j+1}}_{=I} \tilde{H}_j \mathbf{z}\|_2.$$

The minimizer \mathbf{z} of the 2-norm on the right can be computed efficiently, and $\mathbf{x}_j = Q_j \mathbf{z}$.

More details are provided, e.g., in [Mey00].



References I

- [Ber09] Dennis S. Bernstein, *Matrix Mathematics: Theory, Facts, and Formulas*, 2nd ed., Princeton University Press, Princeton, N.J., July 2009.
- [Mey00] Carl D. Meyer, *Matrix Analysis and Applied Linear Algebra*, SIAM, Philadelphia, PA, 2000.
- [MVL78] C. Moler and C. Van Loan, *Nineteen Dubious Ways to Compute the Exponential of a Matrix*, SIAM Rev. **20** (1978), no. 4, 801–836.
- [MVL03] _____, *Nineteen Dubious Ways to Compute the Exponential of a Matrix, Twenty-Five Years Later*, SIAM Rev. **45** (2003), no. 1, 3–49.
- [YM14] Xiyang I. A. Yang and Rajat Mittal, *Acceleration of the Jacobi iterative method by factors exceeding 100 using scheduled relaxation*, Journal of Computational Physics **274** (2014), 695–708.

