

數值分析 MA-3021

Chapter 5. Direct and Iterative Methods for Solving Linear Systems

§5.1 Introduction - Review on Linear Algebra

§5.2 LU decomposition

§5.3 Norms of Vectors and Matrices

§5.4 Iterative Methods

§5.5 Absolute Error, Relative Error and Condition Number

§5.1 Introduction - Review on Linear Algebra

Notation:

- Let A be a $m \times n$ matrix. Then
 - The (i, j) entry of A is denoted by A_{ij} , a_{ij} or $A(i, j)$.
 - The j -th row of A is denoted by $A(j, :)$.
 - The j -th column of A is denoted by $A(:, j)$.
- The $n \times n$ identity matrix is denoted by I_n or $I_{n \times n}$. When the dimension n is clear, $n \times n$ we sometimes also use I to denote the identity matrix.

§5.1 Introduction - Review on Linear Algebra

If A and B are two matrices such that $AB = I$, then we say that B is a right inverse of A and that A is a left inverse of B . For example,

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ \alpha & \beta \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I_{2 \times 2}, \quad \forall \alpha, \beta \in \mathbb{R}.$$

$$\begin{bmatrix} 1 & 0 & \alpha \\ 0 & 1 & \beta \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I_{2 \times 2}, \quad \forall \alpha, \beta \in \mathbb{R}.$$

Notice that right inverse and left inverse may not be unique.

§5.1 Introduction - Review on Linear Algebra

Theorem

A square matrix can possess at most one right inverse.

Proof.

Let $AB = I$. Then $\sum_{j=1}^n b_{jk}A(:,j) = I(:,k)$ for all $1 \leq k \leq n$. So, the columns of A form a basis for \mathbb{R}^n . Therefore, the coefficients b_{jk} above are uniquely determined. \square

Theorem

If A and B are square matrices such that $AB = I$, then $BA = I$.

Proof.

Let $C = BA - I + B$. Then $AC = ABA - AI + AB = A - A + I = I$. Since right inverse for square matrix is at most one, $B = C$. Hence, $C = BA - I + B = BA - I + C$; that is, $BA = I$. \square

§5.1 Introduction - Review on Linear Algebra

Theorem

A square matrix can possess at most one right inverse.

Proof.

Let $AB = I$. Then $\sum_{j=1}^n b_{jk}A(:,j) = I(:,k)$ for all $1 \leq k \leq n$. So, the columns of A form a basis for \mathbb{R}^n . Therefore, the coefficients b_{jk} above are uniquely determined. \square

Theorem

If A and B are square matrices such that $AB = I$, then $BA = I$.

Proof.

Let $C = BA - I + B$. Then $AC = ABA - AI + AB = A - A + I = I$. Since right inverse for square matrix is at most one, $B = C$. Hence, $C = BA - I + B = BA - I + C$; that is, $BA = I$. \square

§5.1 Introduction - Review on Linear Algebra

- 1 If a square matrix A has a right inverse B , then B is unique and $BA = AB = I$. We then call B the inverse of A and say that A is invertible or nonsingular. We denote $B = A^{-1}$.
- 2 If A is invertible, then the system of equations $A\mathbf{x} = \mathbf{b}$ has the solution $\mathbf{x} = A^{-1}\mathbf{b}$. If A^{-1} is not available, then in general, A^{-1} should not be computed solely for the purpose of obtaining x .
- 3 How do we get this A^{-1} ?

§5.1 Introduction - Review on Linear Algebra

- 1 Let two linear systems be given, each consisting of n equations with n unknowns:

$$Ax = b \quad \text{and} \quad Bx = d.$$

If the two systems have precisely the same solutions, we call them equivalent systems.

- 2 Note that A and B can be very different.
- 3 Thus, to solve a linear system of equations, we can instead solve any equivalent system. This simple idea is at the heart of our numerical procedures.

§5.1 Introduction - Review on Linear Algebra

Let \mathcal{E}_i denote the i -th equation in the system $A\mathbf{x} = \mathbf{b}$. The following are the **elementary operations** which can be performed:

- Interchanging two equations in the system: $\mathcal{E}_i \leftrightarrow \mathcal{E}_j$
- Multiplying an equation by a **nonzero** number: $\lambda\mathcal{E}_i \rightarrow \mathcal{E}_i$
- Adding to an equation a multiple of some other equation: $\mathcal{E}_i + \lambda\mathcal{E}_j \rightarrow \mathcal{E}_i$.

Theorem

If one system of equations is obtained from another by a finite sequence of elementary operations, then the two systems are equivalent.

§5.1 Introduction - Review on Linear Algebra

- 1 An **elementary matrix** is defined to be an $n \times n$ matrix that arises when an elementary operation is applied to the $n \times n$ identity matrix.
- 2 The elementary operations expressed in terms of the rows of matrix A are:
 - The interchange of two rows in A : $A(i, :) \leftrightarrow A(j, :)$;
 - Multiplying one row by a **nonzero** constant: $\lambda A(i, :) \rightarrow A(:, i)$;
 - Adding to one row a multiple of another:

$$A(i, :) + \lambda A(j, :) \rightarrow A(i, :).$$

- 3 Each elementary row operation on A can be accomplished by multiplying A on the left by an elementary matrix.

§5.1 Introduction - Review on Linear Algebra

Example

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{31} & a_{32} & a_{33} \\ a_{21} & a_{22} & a_{23} \end{bmatrix}.$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & \lambda & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ \lambda a_{21} & \lambda a_{22} & \lambda a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}.$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & \lambda & 1 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ \lambda a_{21} + a_{31} & \lambda a_{22} + a_{32} & \lambda a_{23} + a_{33} \end{bmatrix}.$$

§5.1 Introduction - Review on Linear Algebra

- 1 If matrix A is invertible, then there exists a sequence of elementary row operations can be applied to A , reducing it to the identity matrix I ,

$$E_m E_{m-1} \cdots E_2 E_1 A = I.$$

- 2 This gives us an equation for computing the inverse of a matrix:

$$A^{-1} = E_m E_{m-1} \cdots E_2 E_1 = E_m E_{m-1} \cdots E_2 E_1 I.$$

Remark: This is not a practical method to compute A^{-1} .

§5.1 Introduction - Review on Linear Algebra

Let $A \in \mathbb{C}^{n \times n}$ be a square matrix. If there exists a nonzero vector $\mathbf{x} \in \mathbb{C}^n$ and a scalar $\lambda \in \mathbb{C}$ such that

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

then λ is called an **eigenvalue** of A and \mathbf{x} is called the corresponding **eigenvector** of A .

Remark: Computing λ and \mathbf{x} is a major task in numerical linear algebra.

§5.1 Introduction - Review on Linear Algebra

For an $n \times n$ real matrix A , the following properties are equivalent:

- 1 The inverse of A exists; that is, A is nonsingular;
- 2 The determinant of A is nonzero;
- 3 The rows of A form a basis for \mathbb{R}^n ;
- 4 The columns of A form a basis for \mathbb{R}^n ;
- 5 As a map from \mathbb{R}^n to \mathbb{R}^n , A is injective (one to one);
- 6 As a map from \mathbb{R}^n to \mathbb{R}^n , A is surjective (onto);
- 7 The equation $A\mathbf{x} = \mathbf{0}$ implies $\mathbf{x} = \mathbf{0}$;
- 8 For each $\mathbf{b} \in \mathbb{R}^n$, there is exactly one $\mathbf{x} \in \mathbb{R}^n$ such that $A\mathbf{x} = \mathbf{b}$;
- 9 A is a product of elementary matrices;
- 10 0 is not an eigenvalue of A .

§5.1 Introduction - Review on Linear Algebra

There are some easy-to-solve systems:

① Diagonal Structure

$$\begin{bmatrix} a_{11} & 0 & 0 & \cdots & 0 \\ 0 & a_{22} & 0 & \cdots & 0 \\ 0 & 0 & a_{33} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ \vdots \\ b_n \end{bmatrix}.$$

The solution is: (provided $a_{ii} \neq 0$ for all $i = 1, 2, \dots, n$)

$$\mathbf{x} = \left(\frac{b_1}{a_{11}}, \frac{b_2}{a_{22}}, \frac{b_3}{a_{33}}, \dots, \frac{b_n}{a_{nn}} \right)^T.$$

- If $a_{ii} = 0$ for some index i , and if $b_i = 0$ also, then x_i can be any real number. The number of solutions is infinity.
- If $a_{ii} = 0$ and $b_i \neq 0$, no solution of the system exists.
- What is the complexity of the method? n divisions.

§5.1 Introduction - Review on Linear Algebra

There are some easy-to-solve systems:

② Lower Triangular Systems

$$\begin{bmatrix} a_{11} & 0 & 0 & \cdots & 0 \\ a_{21} & a_{22} & 0 & \cdots & 0 \\ a_{31} & a_{32} & a_{33} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & a_{n3} & \cdots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ \vdots \\ b_n \end{bmatrix}.$$

Some simple observations:

- If $a_{11} \neq 0$, then we have $x_1 = b_1/a_{11}$.
- Once we have x_1 , we can simplify the second equation, $x_2 = (b_2 - a_{21}x_1)/a_{22}$, provided that $a_{22} \neq 0$.

Similarly, $x_3 = (b_3 - a_{31}x_1 - a_{32}x_2)/a_{33}$, provided that $a_{33} \neq 0$.

In general, to find the solution to this system, we use **forward substitution** (assume that $a_{ii} \neq 0$ for all i).

§5.1 Introduction - Review on Linear Algebra

There are some easy-to-solve systems:

② Lower Triangular Systems (cont'd)

- Algorithm of forward substitution:

input $n, (a_{ij}), b = (b_1, b_2, \dots, b_n)^T$

for $i = 1$ **to** n **do**

$$x_i \leftarrow \left(b_i - \sum_{j=1}^{i-1} a_{ij}x_j \right) / a_{ii}$$

end do

output $x = (x_1, x_2, \dots, x_n)^T$

- Complexity of forward substitution:
 - n divisions.
 - the number of **multiplications**: 0 for x_1 , 1 for x_2 , 2 for x_3 , \dots
total = $0 + 1 + 2 + \dots + (n-1) \approx (n+1)n/2 = \mathcal{O}(n^2)$.
 - the number of **subtractions**: same as the number of multiplications = $\mathcal{O}(n^2)$.

Forward substitution is an $\mathcal{O}(n^2)$ algorithm.

- Remark:** forward substitution is a sequential algorithm (not parallel at all).

§5.1 Introduction - Review on Linear Algebra

There are some easy-to-solve systems:

③ Upper Triangular Systems

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & \cdots & a_{1n} \\ 0 & a_{22} & a_{23} & \cdots & a_{2n} \\ 0 & 0 & a_{33} & \cdots & a_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ \vdots \\ b_n \end{bmatrix}.$$

- The formal algorithm to solve for x is called **backward substitution**. It is also an $\mathcal{O}(n^2)$ algorithm.
- Assume that $a_{ii} \neq 0$ for all i . Algorithm:

input $n, (a_{ij}), b = (b_1, b_2, \dots, b_n)^\top$

for $i = n : -1 : 1$ **do**

$$x_i \leftarrow \left(b_i - \sum_{j=i+1}^n a_{ij}x_j \right) / a_{ii}$$

end do

output $x = (x_1, x_2, \dots, x_n)^\top$

§5.2 LU Decomposition

LU decomposition (factorization):

Suppose that A can be factored into the product of a lower triangular matrix L and an upper triangular matrix U :

$$A = LU.$$

Then, $Ax = LUx = L(Ux)$. Thus, to solve the system of equations $Ax = b$, it is enough to solve this problem in two stages:

$$Lz = b \quad \text{solve for } z,$$

$$Ux = z \quad \text{solve for } x.$$

§5.2 LU Decomposition

Example (Basic Gaussian elimination)

Let $A^{(1)} = (a_{ij}^{(1)}) = A = (a_{ij})$ and $b^{(1)} = b$. Consider the following linear system $Ax = b$:

$$\begin{bmatrix} 6 & -2 & 2 & 4 \\ 12 & -8 & 6 & 10 \\ 3 & -13 & 9 & 3 \\ -6 & 4 & 1 & -18 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 12 \\ 34 \\ 27 \\ -38 \end{bmatrix}.$$

pivot row = row1; pivot element: $a_{11}^{(1)} = 6$.

row2 $- (12/6) \times$ row1 \rightarrow row2.

row3 $- (3/6) \times$ row1 \rightarrow row3.

row4 $- (-6/6) \times$ row1 \rightarrow row4.

$$\Rightarrow \begin{bmatrix} 6 & -2 & 2 & 4 \\ 0 & -4 & 2 & 2 \\ 0 & -12 & 8 & 1 \\ 0 & 2 & 3 & -14 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 12 \\ 10 \\ 21 \\ -26 \end{bmatrix}.$$

multipliers: $12/6$, $3/6$, $(-6)/6$

§5.2 LU Decomposition

Example (Basic Gaussian elimination - cont'd)

We have the following equivalent system $A^{(2)}\mathbf{x} = \mathbf{b}^{(2)}$:

$$\begin{bmatrix} 6 & -2 & 2 & 4 \\ 0 & -4 & 2 & 2 \\ 0 & -12 & 8 & 1 \\ 0 & 2 & 3 & -14 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 12 \\ 10 \\ 21 \\ -26 \end{bmatrix}.$$

pivot row = row2; pivot element $a_{22}^{(2)} = -4$.

row3 $- (-12/-4) \times$ row2 \rightarrow row3.

row4 $- (2/-4) \times$ row2 \rightarrow row4.

$$\Rightarrow \begin{bmatrix} 6 & -2 & 2 & 4 \\ 0 & -4 & 2 & 2 \\ 0 & 0 & 2 & -5 \\ 0 & 0 & 4 & -13 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 12 \\ 10 \\ -9 \\ -21 \end{bmatrix}.$$

multiplier: $(-12)/(-4)$, $2/(-4)$

§5.2 LU Decomposition

Example (Basic Gaussian elimination - cont'd)

We have the following equivalent system $A^{(3)}\mathbf{x} = \mathbf{b}^{(3)}$:

$$\begin{bmatrix} 6 & -2 & 2 & 4 \\ 0 & -4 & 2 & 2 \\ 0 & 0 & 2 & -5 \\ 0 & 0 & 4 & -13 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 12 \\ 10 \\ -9 \\ -21 \end{bmatrix}.$$

pivot row = row3; pivot element $a_{33}^{(3)} = 2$.

row4 $- (4/2) \times$ row3 \rightarrow row4.

$$\Rightarrow \begin{bmatrix} 6 & -2 & 2 & 4 \\ 0 & -4 & 2 & 2 \\ 0 & 0 & 2 & -5 \\ 0 & 0 & 0 & -3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 12 \\ 10 \\ -9 \\ -3 \end{bmatrix}.$$

multiplier: $4/2$

§5.2 LU Decomposition

Example (Basic Gaussian elimination - cont'd)

Finally, we have the following equivalent upper triangular system $A^{(4)}\mathbf{x} = \mathbf{b}^{(4)}$:

$$\begin{bmatrix} 6 & -2 & 2 & 4 \\ 0 & -4 & 2 & 2 \\ 0 & 0 & 2 & -5 \\ 0 & 0 & 0 & -3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 12 \\ 10 \\ -9 \\ -3 \end{bmatrix}.$$

Using the backward substitution, we have

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 1 \\ -3 \\ -2 \\ 1 \end{bmatrix}.$$

§5.2 LU Decomposition

Example (Basic Gaussian elimination - cont'd)

Display the multipliers in an unit lower triangular matrix $L = (\ell_{ij})$:

$$L = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ \frac{1}{2} & 3 & 1 & 0 \\ -1 & -\frac{1}{2} & 2 & 1 \end{bmatrix}.$$

Let $U = (u_{ij})$ be the final upper triangular matrix $A^{(4)}$. Then we have

$$U = \begin{bmatrix} 6 & -2 & 2 & 4 \\ 0 & -4 & 2 & 2 \\ 0 & 0 & 2 & -5 \\ 0 & 0 & 0 & -3 \end{bmatrix}$$

and one can check that $A = LU$ (the Doolittle Decomposition).

§5.2 LU Decomposition

Remark:

- 1 The entire elimination process will break down if any of the pivot elements are 0.
- 2 The total number of arithmetic operations:
 - multiplication and division = $\frac{n^3}{3} - \frac{n}{3} \left(\sum_{k=1}^{n-1} k(k+1) \right)$;
 - addition and subtraction = $\frac{n^3}{3} - \frac{n^2}{2} + \frac{n}{6} \left(\sum_{k=1}^{n-1} k^2 \right)$.

Therefore, the Gauss Elimination is an $O(n^3)$ algorithm.

§5.3 Norms on Vectors and Matrices

Definition

A **normed vector space** $(\mathcal{V}, \|\cdot\|)$ is a vector space \mathcal{V} over field \mathbb{F} associated with a function $\|\cdot\| : \mathcal{V} \rightarrow \mathbb{R}$ such that

- ① $\|\mathbf{x}\| \geq 0$ for all $\mathbf{x} \in \mathcal{V}$.
- ② $\|\mathbf{x}\| = 0$ if and only if $\mathbf{x} = \mathbf{0}$.
- ③ $\|\lambda \cdot \mathbf{x}\| = |\lambda| \cdot \|\mathbf{x}\|$ for all $\lambda \in \mathbb{F}$ and $\mathbf{x} \in \mathcal{V}$.
- ④ $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ for all $\mathbf{x}, \mathbf{y} \in \mathcal{V}$.

A function $\|\cdot\|$ satisfies ①–④ is called a **norm** on \mathcal{V} .

Remark: The norm of a vector can be viewed as the length of that vector. Moreover, the norm induces the concept of distance on the vector space: the distance between two points \mathbf{x} and \mathbf{y} in a normed vector space $(\mathcal{V}, \|\cdot\|)$ is defined by $d(\mathbf{x}, \mathbf{y}) \equiv \|\mathbf{x} - \mathbf{y}\|$.

§5.3 Norms on Vectors and Matrices

Definition

A **normed vector space** $(\mathcal{V}, \|\cdot\|)$ is a vector space \mathcal{V} over field \mathbb{F} associated with a function $\|\cdot\| : \mathcal{V} \rightarrow \mathbb{R}$ such that

- ① $\|\mathbf{x}\| \geq 0$ for all $\mathbf{x} \in \mathcal{V}$.
- ② $\|\mathbf{x}\| = 0$ if and only if $\mathbf{x} = \mathbf{0}$.
- ③ $\|\lambda \cdot \mathbf{x}\| = |\lambda| \cdot \|\mathbf{x}\|$ for all $\lambda \in \mathbb{F}$ and $\mathbf{x} \in \mathcal{V}$.
- ④ $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ for all $\mathbf{x}, \mathbf{y} \in \mathcal{V}$.

A function $\|\cdot\|$ satisfies ①–④ is called a **norm** on \mathcal{V} .

Remark: The norm of a vector can be viewed as the length of that vector. Moreover, the norm induces the concept of distance on the vector space: **the distance between two points \mathbf{x} and \mathbf{y} in a normed vector space $(\mathcal{V}, \|\cdot\|)$ is defined by $d(\mathbf{x}, \mathbf{y}) \equiv \|\mathbf{x} - \mathbf{y}\|$.**

§5.3 Norms on Vectors and Matrices

Example

① Let $\mathbf{x} = (x_1, x_2, \dots, x_n)^\top \in \mathbb{R}^n$:

- The 2-norm (Euclidean norm, or ℓ^2 norm): $\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$
- The infinity norm (ℓ^∞ -norm): $\|\mathbf{x}\|_\infty = \max_{1 \leq i \leq n} |x_i|$
- The 1-norm (ℓ^1 -norm): $\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$
- The p -norm (ℓ^p -norm), $1 \leq p < \infty$, is $\|\mathbf{x}\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}}$

② Let $\mathbf{x} = (x_1, x_2, \dots, x_n)^\top$, $\mathbf{y} = (y_1, y_2, \dots, y_n)^\top \in \mathbb{R}^n$. Then

- $\|\mathbf{x} - \mathbf{y}\|_2 = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$
- $\|\mathbf{x} - \mathbf{y}\|_\infty = \max_{1 \leq i \leq n} |x_i - y_i|$
- $\|\mathbf{x} - \mathbf{y}\|_1 = \sum_{i=1}^n |x_i - y_i|$
- $\|\mathbf{x} - \mathbf{y}\|_p = \left(\sum_{i=1}^n |x_i - y_i|^p \right)^{\frac{1}{p}}$

§5.3 Norms on Vectors and Matrices

Example

$(\mathbb{R}^2, \|\cdot\|_p)$ is a normed vector space. Consider the ball centered at $\mathbf{x}_0 = \mathbf{0}$ with radius 1 and $p = 1$, $p = 2$ and $p = \infty$ respectively.

- $p = 1$: $\|\mathbf{x} - \mathbf{x}_0\|_1 = |x_1| + |x_2|$.
- $p = 2$: $\|\mathbf{x} - \mathbf{x}_0\|_2 = \sqrt{x_1^2 + x_2^2}$.
- $p = \infty$: $\|\mathbf{x} - \mathbf{x}_0\|_\infty = \max\{|x_1|, |x_2|\}$.

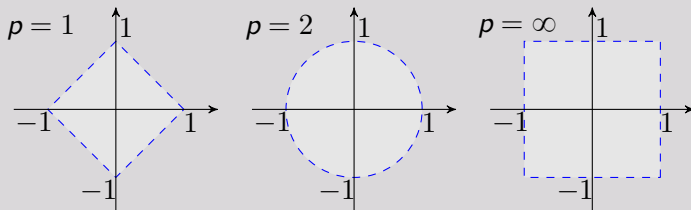


Figure 1: The 1-ball about 0 in \mathbb{R}^2 with different p

§5.3 Norms on Vectors and Matrices

Example

Let A be an invertible $n \times n$ matrix. For a given norm $\|\cdot\|_{\mathbb{R}^n}$ on \mathbb{R}^n , define a map $\|\|\cdot\|\| : \mathbb{R}^n \rightarrow \mathbb{R}$ by

$$\|\|\mathbf{x}\|\| = \|A\mathbf{x}\|_{\mathbb{R}^n}.$$

Then $\|\|\cdot\|\|$ is a norm on \mathbb{R}^n .

Definition

Let $(\mathcal{V}, \|\cdot\|)$ be a normed vector space, and $\{\mathbf{x}^{(n)}\}_{n=1}^{\infty}$ be a sequence in \mathcal{V} . Then $\{\mathbf{x}^{(n)}\}_{n=1}^{\infty}$ is said to converge to a vector $\mathbf{x} \in \mathcal{V}$, denoted by $\lim_{n \rightarrow \infty} \mathbf{x}^{(n)} = \mathbf{x}$, if for every $\varepsilon > 0$, there exists $N > 0$ such that

$$\|\mathbf{x}^{(n)} - \mathbf{x}\| < \varepsilon \quad \text{whenever } n \geq N.$$

Sequence $\{\mathbf{x}^{(n)}\}_{n=1}^{\infty}$ in \mathcal{V} is said to be convergent if there exists $\mathbf{x} \in \mathcal{V}$ such that $\lim_{n \rightarrow \infty} \mathbf{x}^{(n)} = \mathbf{x}$.

§5.3 Norms on Vectors and Matrices

Example

Let A be an invertible $n \times n$ matrix. For a given norm $\|\cdot\|_{\mathbb{R}^n}$ on \mathbb{R}^n , define a map $\|\|\cdot\|\| : \mathbb{R}^n \rightarrow \mathbb{R}$ by

$$\|\|\mathbf{x}\|\| = \|A\mathbf{x}\|_{\mathbb{R}^n}.$$

Then $\|\|\cdot\|\|$ is a norm on \mathbb{R}^n .

Definition

Let $(\mathcal{V}, \|\cdot\|)$ be a normed vector space, and $\{\mathbf{x}^{(n)}\}_{n=1}^{\infty}$ be a sequence in \mathcal{V} . Then $\{\mathbf{x}^{(n)}\}_{n=1}^{\infty}$ is said to converge to a vector $\mathbf{x} \in \mathcal{V}$, denoted by $\lim_{n \rightarrow \infty} \mathbf{x}^{(n)} = \mathbf{x}$, if for every $\varepsilon > 0$, there exists $N > 0$ such that

$$\|\mathbf{x}^{(n)} - \mathbf{x}\| < \varepsilon \quad \text{whenever } n \geq N.$$

Sequence $\{\mathbf{x}^{(n)}\}_{n=1}^{\infty}$ in \mathcal{V} is said to be convergent if there exists $\mathbf{x} \in \mathcal{V}$ such that $\lim_{n \rightarrow \infty} \mathbf{x}^{(n)} = \mathbf{x}$.

§5.3 Norms on Vectors and Matrices

Definition

Let $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$, $(\mathcal{W}, \|\cdot\|_{\mathcal{W}})$ be normed vector spaces, $A \subseteq \mathcal{V}$, and $f: A \rightarrow \mathcal{W}$ be a \mathcal{W} -valued function. f is said to be continuous at $\mathbf{a} \in A$ if for every $\varepsilon > 0$ there exists $\delta > 0$ such that

$$\|f(\mathbf{x}) - f(\mathbf{a})\|_{\mathcal{W}} < \varepsilon \quad \text{whenever} \quad \|\mathbf{x} - \mathbf{a}\|_{\mathcal{V}} < \delta \quad \text{and} \quad \mathbf{x} \in A.$$

Definition

Two norms $\|\cdot\|$ and $\|\!\|\!\cdot\!\!\|$ on a vector space \mathcal{V} are called equivalent if there are positive constants C_1 and C_2 such that

$$C_1\|\mathbf{x}\| \leq \|\!\|\!\mathbf{x}\!\!\| \leq C_2\|\mathbf{x}\| \quad \forall \mathbf{x} \in \mathcal{V}.$$

Remark: Equivalent norms induce the same concept of convergence of sequences, continuity of functions, and so on. For example, if $\{x^{(k)}\}_{k=1}^{\infty}$ is a convergent sequence in $(\mathcal{V}, \|\cdot\|_1)$ and $\|\cdot\|_2$ is an equivalent norm of $\|\cdot\|_1$, then $\{x^{(k)}\}_{k=1}^{\infty}$ is convergent in $(\mathcal{V}, \|\cdot\|_2)$.

§5.3 Norms on Vectors and Matrices

Definition

Let $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$, $(\mathcal{W}, \|\cdot\|_{\mathcal{W}})$ be normed vector spaces, $A \subseteq \mathcal{V}$, and $f: A \rightarrow \mathcal{W}$ be a \mathcal{W} -valued function. f is said to be continuous at $\mathbf{a} \in A$ if for every $\varepsilon > 0$ there exists $\delta > 0$ such that

$$\|f(\mathbf{x}) - f(\mathbf{a})\|_{\mathcal{W}} < \varepsilon \quad \text{whenever} \quad \|\mathbf{x} - \mathbf{a}\|_{\mathcal{V}} < \delta \quad \text{and} \quad \mathbf{x} \in A.$$

Definition

Two norms $\|\cdot\|$ and $\|\!\|\!\cdot\!\!\|$ on a vector space \mathcal{V} are called equivalent if there are positive constants C_1 and C_2 such that

$$C_1\|\mathbf{x}\| \leq \|\!\|\!\mathbf{x}\!\!\| \leq C_2\|\mathbf{x}\| \quad \forall \mathbf{x} \in \mathcal{V}.$$

Remark: Equivalent norms induce the same concept of convergence of sequences, continuity of functions, and so on. For example, if $\{x^{(k)}\}_{k=1}^{\infty}$ is a convergent sequence in $(\mathcal{V}, \|\cdot\|_1)$ and $\|\cdot\|_2$ is an equivalent norm of $\|\cdot\|_1$, then $\{x^{(k)}\}_{k=1}^{\infty}$ is convergent in $(\mathcal{V}, \|\cdot\|_2)$.

§5.3 Norms on Vectors and Matrices

Definition

Let $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$, $(\mathcal{W}, \|\cdot\|_{\mathcal{W}})$ be normed vector spaces, $A \subseteq \mathcal{V}$, and $f: A \rightarrow \mathcal{W}$ be a \mathcal{W} -valued function. f is said to be continuous at $\mathbf{a} \in A$ if for every $\varepsilon > 0$ there exists $\delta > 0$ such that

$$\|f(\mathbf{x}) - f(\mathbf{a})\|_{\mathcal{W}} < \varepsilon \quad \text{whenever} \quad \|\mathbf{x} - \mathbf{a}\|_{\mathcal{V}} < \delta \quad \text{and} \quad \mathbf{x} \in A.$$

Definition

Two norms $\|\cdot\|$ and $\|\!\|\cdot\!\|$ on a vector space \mathcal{V} are called equivalent if there are positive constants C_1 and C_2 such that

$$C_1\|\mathbf{x}\| \leq \|\!\|\mathbf{x}\!\| \leq C_2\|\mathbf{x}\| \quad \forall \mathbf{x} \in \mathcal{V}.$$

Remark: Equivalent norms induce the same concept of convergence of sequences, continuity of functions, and so on. For example, if $\{x^{(k)}\}_{k=1}^{\infty}$ is a convergent sequence in $(\mathcal{V}, \|\cdot\|_1)$ and $\|\cdot\|_2$ is an equivalent norm of $\|\cdot\|_1$, then $\{x^{(k)}\}_{k=1}^{\infty}$ is convergent in $(\mathcal{V}, \|\cdot\|_2)$.

§5.3 Norms on Vectors and Matrices

Theorem

Any two norms on a finite dimensional real (or complex) normed vector space \mathcal{V} are equivalent.

Proof.

Let $\{\mathbf{e}_k\}_{k=1}^N$ be a basis of \mathcal{V} . For each $\mathbf{x} \in \mathcal{V}$, we write $\mathbf{x} = \sum_{k=1}^N x_k \mathbf{e}_k$ and define a function $\|\cdot\|_2 : \mathcal{V} \rightarrow \mathbb{R}$ by

$$\|\mathbf{x}\|_2 = \left(\sum_{k=1}^N |x_k|^2 \right)^{\frac{1}{2}}.$$

Then

- ① $\|\mathbf{x}\|_2 \geq 0$ for all $\mathbf{x} \in \mathcal{V}$, and $\|\mathbf{x}\|_2 = 0$ if and only if $\mathbf{x} = \mathbf{0}$.
- ② $\|\lambda \mathbf{x}\|_2 = |\lambda| \|\mathbf{x}\|_2$ for all $\lambda \in \mathbb{R}$ (or \mathbb{C}) and $\mathbf{x} \in \mathcal{V}$.
- ③ $\|\mathbf{x} + \mathbf{y}\|_2 \leq \|\mathbf{x}\|_2 + \|\mathbf{y}\|_2$ for all $\mathbf{x}, \mathbf{y} \in \mathcal{V}$ because of the Cauchy-Schwarz inequality. □

§5.3 Norms on Vectors and Matrices

Theorem

Any two norms on a finite dimensional real (or complex) normed vector space \mathcal{V} are equivalent.

Proof (cont'd).

Therefore, $\|\cdot\|_2$ is a norm on \mathcal{V} . It then suffices to show that any norm $\|\cdot\|$ on \mathcal{V} is equivalent to $\|\cdot\|_2$:

$$\text{if } C_1\|\mathbf{x}\| \leq \|\mathbf{x}\|_2 \leq C_2\|\mathbf{x}\| \quad \text{and} \quad C_3\|\mathbf{x}\| \leq \|\mathbf{x}\|_2 \leq C_4\|\mathbf{x}\|,$$

$$\text{then } \frac{C_1}{C_4}\|\mathbf{x}\| \leq \|\mathbf{x}\| \leq \frac{C_2}{C_3}\|\mathbf{x}\|.$$

By the definition of norms and the Cauchy-Schwarz inequality,

$$\|\mathbf{x}\| \leq \sum_{k=1}^N |x_k| \|\mathbf{e}_k\| \leq \|\mathbf{x}\|_2 \left(\sum_{k=1}^N \|\mathbf{e}_k\|^2 \right)^{\frac{1}{2}};$$

thus letting $C_2 = \left(\sum_{k=1}^N \|\mathbf{e}_k\|^2 \right)^{\frac{1}{2}}$ we have $\|\mathbf{x}\| \leq C_2\|\mathbf{x}\|_2$. □

§5.3 Norms on Vectors and Matrices

Theorem

Any two norms on a finite dimensional real (or complex) normed vector space \mathcal{V} are equivalent.

Proof (cont'd).

Therefore, $\|\cdot\|_2$ is a norm on \mathcal{V} . It then suffices to show that any norm $\|\cdot\|$ on \mathcal{V} is equivalent to $\|\cdot\|_2$:

$$\text{if } C_1\|\mathbf{x}\| \leq \|\mathbf{x}\|_2 \leq C_2\|\mathbf{x}\| \quad \text{and} \quad C_3\|\mathbf{x}\| \leq \|\mathbf{x}\|_2 \leq C_4\|\mathbf{x}\|,$$

$$\text{then } \frac{C_1}{C_4}\|\mathbf{x}\| \leq \|\mathbf{x}\| \leq \frac{C_2}{C_3}\|\mathbf{x}\|.$$

By the definition of norms and the Cauchy-Schwarz inequality,

$$\|\mathbf{x}\| \leq \sum_{k=1}^N |x_k| \|\mathbf{e}_k\| \leq \|\mathbf{x}\|_2 \left(\sum_{k=1}^N \|\mathbf{e}_k\|^2 \right)^{\frac{1}{2}};$$

thus letting $C_2 = \left(\sum_{k=1}^N \|\mathbf{e}_k\|^2 \right)^{\frac{1}{2}}$ we have $\|\mathbf{x}\| \leq C_2\|\mathbf{x}\|_2$. □

§5.3 Norms on Vectors and Matrices

Theorem

Any two norms on a finite dimensional real (or complex) normed vector space \mathcal{V} are equivalent.

Proof (cont'd).

Define $f: (\mathcal{V}, \|\cdot\|_2) \rightarrow \mathbb{R}$ by $f(\mathbf{x}) = \|\mathbf{x}\|$. Then

$$|f(\mathbf{x}) - f(\mathbf{y})| = \left| \|\mathbf{x}\| - \|\mathbf{y}\| \right| \leq \|\mathbf{x} - \mathbf{y}\| \leq C_2 \|\mathbf{x} - \mathbf{y}\|_2$$

which implies that f is continuous. Let \mathbb{S}^{n-1} be the unit sphere $\{\mathbf{x} \in \mathcal{V} \mid \|\mathbf{x}\|_2 = 1\}$. Then \mathbb{S}^{n-1} is (sequentially) compact in $(\mathcal{V}, \|\cdot\|_2)$, so f attains its minimum on \mathbb{S}^{n-1} . Suppose that $\min_{\mathbf{x} \in \mathbb{S}^{n-1}} f(\mathbf{x}) = f(\mathbf{a})$ for some $\mathbf{a} \in \mathbb{S}^{n-1}$. Then $f(\mathbf{a}) > 0$ (for otherwise $\mathbf{a} = \mathbf{0}$), and

$$\left\| \frac{\mathbf{x}}{\|\mathbf{x}\|_2} \right\| = f\left(\frac{\mathbf{x}}{\|\mathbf{x}\|_2}\right) \geq f(\mathbf{a}) \quad \forall \mathbf{x} \in \mathcal{V}$$

which implies that $\|\mathbf{x}\| \geq C_1 \|\mathbf{x}\|_2$ for $C_1 = f(\mathbf{a})$. □

§5.3 Norms on Vectors and Matrices

Theorem

Any two norms on a finite dimensional real (or complex) normed vector space \mathcal{V} are equivalent.

Proof (cont'd).

Define $f: (\mathcal{V}, \|\cdot\|_2) \rightarrow \mathbb{R}$ by $f(\mathbf{x}) = \|\mathbf{x}\|$. Then

$$|f(\mathbf{x}) - f(\mathbf{y})| = \left| \|\mathbf{x}\| - \|\mathbf{y}\| \right| \leq \|\mathbf{x} - \mathbf{y}\| \leq C_2 \|\mathbf{x} - \mathbf{y}\|_2$$

which implies that f is continuous. Let \mathbb{S}^{n-1} be the unit sphere $\{\mathbf{x} \in \mathcal{V} \mid \|\mathbf{x}\|_2 = 1\}$. Then \mathbb{S}^{n-1} is (sequentially) compact in $(\mathcal{V}, \|\cdot\|_2)$, so f attains its minimum on \mathbb{S}^{n-1} . Suppose that $\min_{\mathbf{x} \in \mathbb{S}^{n-1}} f(\mathbf{x}) = f(\mathbf{a})$ for some $\mathbf{a} \in \mathbb{S}^{n-1}$. Then $f(\mathbf{a}) > 0$ (for otherwise $\mathbf{a} = \mathbf{0}$), and

$$\left\| \frac{\mathbf{x}}{\|\mathbf{x}\|_2} \right\| = f\left(\frac{\mathbf{x}}{\|\mathbf{x}\|_2}\right) \geq f(\mathbf{a}) \quad \forall \mathbf{x} \in \mathcal{V}$$

which implies that $\|\mathbf{x}\| \geq C_1 \|\mathbf{x}\|_2$ for $C_1 = f(\mathbf{a})$. □

§5.3 Norms on Vectors and Matrices

Theorem

Any two norms on a finite dimensional real (or complex) normed vector space \mathcal{V} are equivalent.

Proof (cont'd).

Define $f: (\mathcal{V}, \|\cdot\|_2) \rightarrow \mathbb{R}$ by $f(\mathbf{x}) = \|\mathbf{x}\|$. Then

$$|f(\mathbf{x}) - f(\mathbf{y})| = \left| \|\mathbf{x}\| - \|\mathbf{y}\| \right| \leq \|\mathbf{x} - \mathbf{y}\| \leq C_2 \|\mathbf{x} - \mathbf{y}\|_2$$

which implies that f is continuous. Let \mathbb{S}^{n-1} be the unit sphere $\{\mathbf{x} \in \mathcal{V} \mid \|\mathbf{x}\|_2 = 1\}$. Then \mathbb{S}^{n-1} is (sequentially) compact in $(\mathcal{V}, \|\cdot\|_2)$, so f attains its minimum on \mathbb{S}^{n-1} . Suppose that $\min_{\mathbf{x} \in \mathbb{S}^{n-1}} f(\mathbf{x}) = f(\mathbf{a})$ for some $\mathbf{a} \in \mathbb{S}^{n-1}$. Then $f(\mathbf{a}) > 0$ (for otherwise $\mathbf{a} = \mathbf{0}$), and

$$\left\| \frac{\mathbf{x}}{\|\mathbf{x}\|_2} \right\| = f\left(\frac{\mathbf{x}}{\|\mathbf{x}\|_2}\right) \geq f(\mathbf{a}) \quad \forall \mathbf{x} \in \mathcal{V}$$

which implies that $\|\mathbf{x}\| \geq C_1 \|\mathbf{x}\|_2$ for $C_1 = f(\mathbf{a})$. □

§5.3 Norms on Vectors and Matrices

Theorem

Any two norms on a finite dimensional real (or complex) normed vector space \mathcal{V} are equivalent.

Proof (cont'd).

Define $f: (\mathcal{V}, \|\cdot\|_2) \rightarrow \mathbb{R}$ by $f(\mathbf{x}) = \|\mathbf{x}\|$. Then

$$|f(\mathbf{x}) - f(\mathbf{y})| = \left| \|\mathbf{x}\| - \|\mathbf{y}\| \right| \leq \|\mathbf{x} - \mathbf{y}\| \leq C_2 \|\mathbf{x} - \mathbf{y}\|_2$$

which implies that f is continuous. Let \mathbb{S}^{n-1} be the unit sphere $\{\mathbf{x} \in \mathcal{V} \mid \|\mathbf{x}\|_2 = 1\}$. Then \mathbb{S}^{n-1} is (sequentially) compact in $(\mathcal{V}, \|\cdot\|_2)$, so f attains its minimum on \mathbb{S}^{n-1} . Suppose that $\min_{\mathbf{x} \in \mathbb{S}^{n-1}} f(\mathbf{x}) = f(\mathbf{a})$ for some $\mathbf{a} \in \mathbb{S}^{n-1}$. Then $f(\mathbf{a}) > 0$ (for otherwise $\mathbf{a} = \mathbf{0}$), and

$$\left\| \frac{\mathbf{x}}{\|\mathbf{x}\|_2} \right\| = f\left(\frac{\mathbf{x}}{\|\mathbf{x}\|_2}\right) \geq f(\mathbf{a}) \quad \forall \mathbf{x} \in \mathcal{V}$$

which implies that $\|\mathbf{x}\| \geq C_1 \|\mathbf{x}\|_2$ for $C_1 = f(\mathbf{a})$. □

§5.3 Norms on Vectors and Matrices

Theorem

Let $\|\cdot\|_{\mathbb{R}^n}$ be a norm on \mathbb{R}^n and $\|\cdot\|_{\mathbb{R}^m}$ be a norm on \mathbb{R}^m . Then

$$\|A\|_{\mathbb{R}^n, \mathbb{R}^m} \equiv \max \{ \|A\mathbf{x}\|_{\mathbb{R}^m} : \mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|_{\mathbb{R}^n} = 1 \}$$

defines a norm on the vector space of all $m \times n$ real matrices.

Proof.

- ① Clearly $\|A\|_{\mathbb{R}^n, \mathbb{R}^m} \geq 0$, and $\|A\| = 0$ if and only if $A = 0$.
- ②
$$\begin{aligned} \|\lambda A\|_{\mathbb{R}^n, \mathbb{R}^m} &= \max \{ \|\lambda A\mathbf{x}\|_{\mathbb{R}^m} : \|\mathbf{x}\|_{\mathbb{R}^n} = 1 \} \\ &= \max \{ |\lambda| \|A\mathbf{x}\|_{\mathbb{R}^m} : \|\mathbf{x}\|_{\mathbb{R}^n} = 1 \} \\ &= |\lambda| \max \{ \|A\mathbf{x}\|_{\mathbb{R}^m} : \|\mathbf{x}\|_{\mathbb{R}^n} = 1 \} = |\lambda| \|A\|_{\mathbb{R}^n, \mathbb{R}^m}. \end{aligned}$$
- ③
$$\begin{aligned} \|A + B\|_{\mathbb{R}^n, \mathbb{R}^m} &= \max \{ \|(A + B)\mathbf{x}\|_{\mathbb{R}^m} : \|\mathbf{x}\|_{\mathbb{R}^n} = 1 \} \\ &\leq \max \{ \|A\mathbf{x}\|_{\mathbb{R}^m} + \|B\mathbf{x}\|_{\mathbb{R}^m} : \|\mathbf{x}\|_{\mathbb{R}^n} = 1 \} \\ &\leq \max \{ \|A\mathbf{x}\|_{\mathbb{R}^m} : \|\mathbf{x}\|_{\mathbb{R}^n} = 1 \} + \max \{ \|B\mathbf{x}\|_{\mathbb{R}^m} : \|\mathbf{x}\|_{\mathbb{R}^n} = 1 \} \\ &= \|A\|_{\mathbb{R}^n, \mathbb{R}^m} + \|B\|_{\mathbb{R}^n, \mathbb{R}^m}. \quad \square \end{aligned}$$

§5.3 Norms on Vectors and Matrices

Remark:

- ① $\|\cdot\|_{\mathbb{R}^n, \mathbb{R}^m}$ is called the matrix norm induced by vector norms $\|\cdot\|_{\mathbb{R}^n}$ and $\|\cdot\|_{\mathbb{R}^m}$. Moreover,

$$\|A\|_{\mathbb{R}^n, \mathbb{R}^m} \equiv \max \{ \|Ax\|_{\mathbb{R}^m} : x \in \mathbb{R}^n, \|x\|_{\mathbb{R}^n} = 1 \}$$

$$\Leftrightarrow \|A\|_{\mathbb{R}^n, \mathbb{R}^m} \equiv \max \left\{ \frac{\|Ax\|_{\mathbb{R}^m}}{\|x\|_{\mathbb{R}^n}} : x \in \mathbb{R}^n, x \neq \mathbf{0} \right\}$$

- ② If $\|\cdot\|_{\mathbb{R}^n} = \|\cdot\|_p$ and $\|\cdot\|_{\mathbb{R}^m} = \|\cdot\|_q$, then $\|A\|_{\mathbb{R}^n, \mathbb{R}^m}$ is simply denoted by $\|A\|_{p,q}$. If in addition $p = q$, then $\|A\|_{p,q}$ is simply denoted by $\|A\|_p$.

Theorem (Additional properties of matrix norms)

Let A be a $m \times n$ matrix, and B be a $n \times k$ matrix.

- ① $\|Ax\|_{\mathbb{R}^m} \leq \|A\|_{\mathbb{R}^n, \mathbb{R}^m} \|x\|_{\mathbb{R}^n}$ for all $x \in \mathbb{R}^n$ (sub-ordinance)
- ② $\|AB\|_{\mathbb{R}^k, \mathbb{R}^m} \leq \|A\|_{\mathbb{R}^n, \mathbb{R}^m} \|B\|_{\mathbb{R}^k, \mathbb{R}^n}$ (sub-multiplicativity)
- ③ $\|I_{n \times n}\|_p = 1$ for all $p \in [1, \infty]$.

§5.3 Norms on Vectors and Matrices

Remark:

- ① $\|\cdot\|_{\mathbb{R}^n, \mathbb{R}^m}$ is called the matrix norm induced by vector norms $\|\cdot\|_{\mathbb{R}^n}$ and $\|\cdot\|_{\mathbb{R}^m}$. Moreover,

$$\|A\|_{\mathbb{R}^n, \mathbb{R}^m} \equiv \max \{ \|Ax\|_{\mathbb{R}^m} : x \in \mathbb{R}^n, \|x\|_{\mathbb{R}^n} = 1 \}$$

$$\Leftrightarrow \|A\|_{\mathbb{R}^n, \mathbb{R}^m} \equiv \max \left\{ \frac{\|Ax\|_{\mathbb{R}^m}}{\|x\|_{\mathbb{R}^n}} : x \in \mathbb{R}^n, x \neq \mathbf{0} \right\}$$

- ② If $\|\cdot\|_{\mathbb{R}^n} = \|\cdot\|_p$ and $\|\cdot\|_{\mathbb{R}^m} = \|\cdot\|_q$, then $\|A\|_{\mathbb{R}^n, \mathbb{R}^m}$ is simply denoted by $\|A\|_{p,q}$. If in addition $p = q$, then $\|A\|_{p,q}$ is simply denoted by $\|A\|_p$.

Theorem (Additional properties of matrix norms)

Let A be a $m \times n$ matrix, and B be a $n \times k$ matrix.

- ① $\|Ax\|_{\mathbb{R}^m} \leq \|A\|_{\mathbb{R}^n, \mathbb{R}^m} \|x\|_{\mathbb{R}^n}$ for all $x \in \mathbb{R}^n$ (sub-ordinance)
- ② $\|AB\|_{\mathbb{R}^k, \mathbb{R}^m} \leq \|A\|_{\mathbb{R}^n, \mathbb{R}^m} \|B\|_{\mathbb{R}^k, \mathbb{R}^n}$ (sub-multiplicativity)
- ③ $\|I_{n \times n}\|_p = 1$ for all $p \in [1, \infty]$.

§5.3 Norms on Vectors and Matrices

Remark:

- ① $\|\cdot\|_{\mathbb{R}^n, \mathbb{R}^m}$ is called the matrix norm induced by vector norms $\|\cdot\|_{\mathbb{R}^n}$ and $\|\cdot\|_{\mathbb{R}^m}$. Moreover,

$$\|A\|_{\mathbb{R}^n, \mathbb{R}^m} \equiv \max \{ \|Ax\|_{\mathbb{R}^m} : x \in \mathbb{R}^n, \|x\|_{\mathbb{R}^n} = 1 \}$$

$$\Leftrightarrow \|A\|_{\mathbb{R}^n, \mathbb{R}^m} \equiv \max \left\{ \frac{\|Ax\|_{\mathbb{R}^m}}{\|x\|_{\mathbb{R}^n}} : x \in \mathbb{R}^n, x \neq \mathbf{0} \right\}$$

- ② If $\|\cdot\|_{\mathbb{R}^n} = \|\cdot\|_p$ and $\|\cdot\|_{\mathbb{R}^m} = \|\cdot\|_q$, then $\|A\|_{\mathbb{R}^n, \mathbb{R}^m}$ is simply denoted by $\|A\|_{p,q}$. If in addition $p = q$, then $\|A\|_{p,q}$ is simply denoted by $\|A\|_p$.

Theorem (Additional properties of matrix norms)

Let A be a $m \times n$ matrix, and B be a $n \times k$ matrix.

- ① $\|Ax\|_{\mathbb{R}^m} \leq \|A\|_{\mathbb{R}^n, \mathbb{R}^m} \|x\|_{\mathbb{R}^n}$ for all $x \in \mathbb{R}^n$ ($\|Ax\| \leq \|A\| \|x\|$)
- ② $\|AB\|_{\mathbb{R}^k, \mathbb{R}^m} \leq \|A\|_{\mathbb{R}^n, \mathbb{R}^m} \|B\|_{\mathbb{R}^k, \mathbb{R}^n}$ ($\|AB\| \leq \|A\| \|B\|$)
- ③ $\|I_{n \times n}\|_p = 1$ for all $p \in [1, \infty]$.

§5.3 Norms on Vectors and Matrices

Example ($\|A\|_\infty$)

Let $A = [a_{ij}]_{m \times n}$ and $\max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}| = \sum_{j=1}^n |a_{kj}|$ for some $1 \leq k \leq m$.

① Let $\mathbf{x} = (\text{sgn}(a_{k1}), \text{sgn}(a_{k2}), \dots, \text{sgn}(a_{kn}))$. Then $\|\mathbf{x}\|_\infty = 1$, and $\|A\mathbf{x}\|_\infty = \sum_{j=1}^n |a_{kj}|$.

② Let $\mathbf{x} = (x_1, \dots, x_n)$. If $\|\mathbf{x}\|_\infty = 1$, then $|x_j| \leq 1$ for all $1 \leq j \leq n$; thus

$$|a_{i1}x_1 + a_{i2}x_2 + \dots + a_{in}x_n| \leq \sum_{j=1}^n |a_{ij}| \leq \sum_{j=1}^n |a_{kj}|.$$

By the definition of matrix norms, ① implies that $\|A\|_\infty \geq \sum_{j=1}^n |a_{kj}|$ while ② implies that $\|A\|_\infty \leq \sum_{j=1}^n |a_{kj}|$. Therefore,

$$\|A\|_\infty = \max \left\{ \sum_{j=1}^n |a_{1j}|, \sum_{j=1}^n |a_{2j}|, \dots, \sum_{j=1}^n |a_{nj}| \right\};$$

that is, $\|A\|_\infty$ is the largest sum of the absolute value of row entries.

§5.3 Norms on Vectors and Matrices

Theorem

For each $\mathbf{x} \in \mathbb{R}^n$,

$$\|\mathbf{x}\|_1 = \max \{ \mathbf{y}^T \mathbf{x} : \|\mathbf{y}\|_\infty = 1 \}, \quad \|\mathbf{x}\|_\infty = \max \{ \mathbf{y}^T \mathbf{x} : \|\mathbf{y}\|_1 = 1 \}.$$

Example ($\|A\|_1$)

By the theorem above,

$$\begin{aligned} \|A\|_1 &= \max_{\|\mathbf{x}\|_1=1} \|A\mathbf{x}\|_1 = \max_{\|\mathbf{x}\|_1=1} \max_{\|\mathbf{y}\|_\infty=1} \mathbf{y}^T A\mathbf{x} \\ &= \max_{\|\mathbf{y}\|_\infty=1} \max_{\|\mathbf{x}\|_1=1} \mathbf{y}^T A\mathbf{x} = \max_{\|\mathbf{y}\|_\infty=1} \max_{\|\mathbf{x}\|_1=1} \mathbf{x}^T A^T \mathbf{y} \\ &= \max_{\|\mathbf{y}\|_\infty=1} \|A^T \mathbf{y}\|_\infty = \|A^T\|_\infty; \end{aligned}$$

thus

$$\|A\|_1 = \max \left\{ \sum_{i=1}^m |a_{i1}|, \sum_{i=1}^m |a_{i2}|, \dots, \sum_{i=1}^m |a_{in}| \right\};$$

that is, $\|A\|_1$ is the largest sum of the absolute value of column entries.

§5.3 Norms on Vectors and Matrices

Example ($\|A\|_2$)

Let A be an $m \times n$ matrix. Then by the definition of the 2-norm,

$$\|A\|_2^2 = \max \{ \|Ax\|_2^2 : \|x\|_2 = 1 \} = \max \{ x^T A^T A x : \|x\|_2 = 1 \}.$$

Since $A^T A$ is an $n \times n$ symmetric matrix, $A^T A$ has n real eigenvalues $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ and corresponding **orthogonal unit** eigenvectors v_1, v_2, \dots, v_n . Then each $x \in \mathbb{R}^n$ can be expressed as

$$x = x_1 v_1 + x_2 v_2 + \dots + x_n v_n \quad (*)$$

and the condition $\|x\|_2 = 1$ is translated into $\sum_{i=1}^n x_i^2 = 1$. Using (*),

$$x^T A^T A x = \lambda_1 x_1^2 + \lambda_2 x_2^2 + \dots + \lambda_n x_n^2$$

whose maximum, under the constraint $\sum_{i=1}^n x_i^2 = 1$, is λ_n . Therefore,

$\|A\|_2$ = the square root of the maximum eigenvalue of $A^T A$.

§5.3 Norms on Vectors and Matrices

Example ($\|A\|_2$)

Let A be an $m \times n$ matrix. Then by the definition of the 2-norm,

$$\|A\|_2^2 = \max \{ \|Ax\|_2^2 : \|x\|_2 = 1 \} = \max \{ x^T A^T A x : \|x\|_2 = 1 \}.$$

Since $A^T A$ is an $n \times n$ symmetric matrix, $A^T A$ has n real eigenvalues $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ and corresponding **orthogonal unit** eigenvectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$. Then each $\mathbf{x} \in \mathbb{R}^n$ can be expressed as

$$\mathbf{x} = x_1 \mathbf{v}_1 + x_2 \mathbf{v}_2 + \dots + x_n \mathbf{v}_n \quad (*)$$

and the condition $\|x\|_2 = 1$ is translated into $\sum_{i=1}^n x_i^2 = 1$. Using (*),

$$x^T A^T A x = \lambda_1 x_1^2 + \lambda_2 x_2^2 + \dots + \lambda_n x_n^2$$

whose maximum, under the constraint $\sum_{i=1}^n x_i^2 = 1$, is λ_n . Therefore,

$\|A\|_2$ = the square root of the maximum eigenvalue of $A^T A$.

§5.3 Norms on Vectors and Matrices

Example ($\|A\|_2$)

Let A be an $m \times n$ matrix. Then by the definition of the 2-norm,

$$\|A\|_2^2 = \max \{ \|Ax\|_2^2 : \|x\|_2 = 1 \} = \max \{ x^T A^T A x : \|x\|_2 = 1 \}.$$

Since $A^T A$ is an $n \times n$ symmetric matrix, $A^T A$ has n real eigenvalues $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ and corresponding **orthogonal unit** eigenvectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$. Then each $\mathbf{x} \in \mathbb{R}^n$ can be expressed as

$$\mathbf{x} = x_1 \mathbf{v}_1 + x_2 \mathbf{v}_2 + \dots + x_n \mathbf{v}_n \quad (*)$$

and the condition $\|x\|_2 = 1$ is translated into $\sum_{i=1}^n x_i^2 = 1$. Using (*),

$$x^T A^T A x = \lambda_1 x_1^2 + \lambda_2 x_2^2 + \dots + \lambda_n x_n^2$$

whose maximum, under the constraint $\sum_{i=1}^n x_i^2 = 1$, is λ_n . Therefore,

$\|A\|_2$ = the square root of the maximum eigenvalue of $A^T A$.

§5.3 Norms on Vectors and Matrices

Example ($\|A\|_2$)

Let A be an $m \times n$ matrix. Then by the definition of the 2-norm,

$$\|A\|_2^2 = \max \{ \|Ax\|_2^2 : \|x\|_2 = 1 \} = \max \{ x^T A^T A x : \|x\|_2 = 1 \}.$$

Since $A^T A$ is an $n \times n$ symmetric matrix, $A^T A$ has n real eigenvalues $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ and corresponding **orthogonal unit** eigenvectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$. Then each $\mathbf{x} \in \mathbb{R}^n$ can be expressed as

$$\mathbf{x} = x_1 \mathbf{v}_1 + x_2 \mathbf{v}_2 + \dots + x_n \mathbf{v}_n \quad (\star)$$

and the condition $\|x\|_2 = 1$ is translated into $\sum_{i=1}^n x_i^2 = 1$. Using (\star) ,

$$\mathbf{x}^T A^T A \mathbf{x} = \lambda_1 x_1^2 + \lambda_2 x_2^2 + \dots + \lambda_n x_n^2$$

whose maximum, under the constraint $\sum_{i=1}^n x_i^2 = 1$, is λ_n . Therefore,

$\|A\|_2$ = the square root of the maximum eigenvalue of $A^T A$.

§5.3 Norms on Vectors and Matrices

Example ($\|A\|_2$)

Let A be an $m \times n$ matrix. Then by the definition of the 2-norm,

$$\|A\|_2^2 = \max \{ \|Ax\|_2^2 : \|x\|_2 = 1 \} = \max \{ x^T A^T A x : \|x\|_2 = 1 \}.$$

Since $A^T A$ is an $n \times n$ symmetric matrix, $A^T A$ has n real eigenvalues $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ and corresponding **orthogonal unit** eigenvectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$. Then each $\mathbf{x} \in \mathbb{R}^n$ can be expressed as

$$\mathbf{x} = x_1 \mathbf{v}_1 + x_2 \mathbf{v}_2 + \dots + x_n \mathbf{v}_n \quad (\star)$$

and the condition $\|x\|_2 = 1$ is translated into $\sum_{i=1}^n x_i^2 = 1$. Using (\star) ,

$$x^T A^T A x = \lambda_1 x_1^2 + \lambda_2 x_2^2 + \dots + \lambda_n x_n^2$$

whose maximum, under the constraint $\sum_{i=1}^n x_i^2 = 1$, is λ_n . Therefore,

$\|A\|_2$ = the square root of the maximum eigenvalue of $A^T A$.

§5.3 Norms on Vectors and Matrices

Example

Consider the matrix $A = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 2 & 1 \\ -1 & 1 & 2 \end{bmatrix}$. Then

$$A^T A = \begin{bmatrix} 1 & 1 & -1 \\ 1 & 2 & 1 \\ 0 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 \\ 1 & 2 & 1 \\ -1 & 1 & 2 \end{bmatrix} = \begin{bmatrix} 3 & 2 & -1 \\ 2 & 6 & 4 \\ -1 & 4 & 5 \end{bmatrix}$$

which implies that the characteristic equation of $A^T A$ is

$$\det(A^T A - \lambda I) = \begin{vmatrix} 3 - \lambda & 2 & -1 \\ 2 & 6 - \lambda & 4 \\ -1 & 4 & 5 - \lambda \end{vmatrix} = -\lambda(\lambda^2 - 14\lambda + 42) = 0.$$

Therefore, the eigenvalues of $A^T A$ are $\lambda = 0, 7 + \sqrt{7}, 7 - \sqrt{7}$; thus

$$\|A\|_2 = \sqrt{\rho(A^T A)} = \sqrt{7 + \sqrt{7}} \approx 3.106.$$

§5.3 Norms on Vectors and Matrices

Example (Frobenius Norm)

Not every norm on the space of $m \times n$ real matrices is of the form $\|\cdot\|_{\mathbb{R}^n, \mathbb{R}^m}$ (called the **natural norm**). For example, the **Frobenius norm**, sometimes also called the **Euclidean norm** (a term unfortunately also used for the vector ℓ^2 -norm), is matrix norm of an $m \times n$ matrix A defined as the square root of the sum of the absolute squares of its elements; that is,

$$\|A\|_F = \left(\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 \right)^{\frac{1}{2}}.$$

This is clear a norm because this is to identify the space of real $m \times n$ matrices as the space \mathbb{R}^{mn} with ℓ^2 -norm. The Frobenius norm can also be computed by

$$\|A\|_F = \sqrt{\text{Tr}(AA^T)},$$

where $\text{Tr}(M)$ is the trace of (a square matrix) M .

§5.3 Norms on Vectors and Matrices

Example (Frobenius Norm)

Not every norm on the space of $m \times n$ real matrices is of the form $\|\cdot\|_{\mathbb{R}^n, \mathbb{R}^m}$ (called the **natural norm**). For example, the **Frobenius norm**, sometimes also called the **Euclidean norm** (a term unfortunately also used for the vector ℓ^2 -norm), is matrix norm of an $m \times n$ matrix A defined as the square root of the sum of the absolute squares of its elements; that is,

$$\|A\|_F = \left(\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 \right)^{\frac{1}{2}}.$$

This is clear a norm because this is to identify the space of real $m \times n$ matrices as the space \mathbb{R}^{mn} with ℓ^2 -norm. The Frobenius norm can also be computed by

$$\|A\|_F = \sqrt{\text{Tr}(AA^T)},$$

where $\text{Tr}(M)$ is the trace of (a square matrix) M .

§5.3 Norms on Vectors and Matrices

Example (Frobenius Norm)

Not every norm on the space of $m \times n$ real matrices is of the form $\|\cdot\|_{\mathbb{R}^n, \mathbb{R}^m}$ (called the **natural norm**). For example, the **Frobenius norm**, sometimes also called the **Euclidean norm** (a term unfortunately also used for the vector ℓ^2 -norm), is matrix norm of an $m \times n$ matrix A defined as the square root of the sum of the absolute squares of its elements; that is,

$$\|A\|_F = \left(\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 \right)^{\frac{1}{2}}.$$

This is clear a norm because this is to **identify the space of real $m \times n$ matrices as the space \mathbb{R}^{mn} with ℓ^2 -norm**. The Frobenius norm can also be computed by

$$\|A\|_F = \sqrt{\text{Tr}(AA^T)},$$

where $\text{Tr}(M)$ is the trace of (a square matrix) M .

§5.3 Norms on Vectors and Matrices

Definition

The spectral radius of a square matrix is the largest absolute value of its eigenvalues. The spectral radius of A is denoted by $\rho(A)$.

Theorem

Let A be an $m \times n$ real matrix. Then $\|A\|_2 = \sqrt{\rho(A^T A)}$.

Remark: The ℓ^2 -matrix norm is also called the spectral norm.

Corollary

If A is a real symmetric matrix, then $\|A\|_2 = \rho(A)$.

Theorem

$\rho(A) \leq \|A\|$ for any real square matrix A and natural norm $\|\cdot\|$.

§5.3 Norms on Vectors and Matrices

Theorem

Let A be a real square matrix. Then for every $\varepsilon > 0$ there exists a (subordinate) matrix norm $\|\cdot\|$ such that $\|A\| \leq \rho(A) + \varepsilon$.

Proof.

Let A be an $n \times n$ real matrix. The Jordan canonical form of A is

$$A = S \begin{bmatrix} J_{n_1}(\lambda_1) & 0 & \cdots & 0 \\ 0 & J_{n_2}(\lambda_2) & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & J_{n_k}(\lambda_k) \end{bmatrix} S^{-1},$$

where S is an invertible matrix, $\lambda_1, \lambda_2, \dots, \lambda_k$ are (complex) eigenvalues of A , $n_1 + n_2 + \cdots + n_k = n$, and $J_{n_j}(\lambda_j)$ are Jordan blocks of size $n_j \times n_j$. □

§5.3 Norms on Vectors and Matrices

Proof (cont'd).

For each $m \in \mathbb{N}$ and $\eta > 0$, define

$$D(\eta) = \begin{bmatrix} D_{n_1}(\eta) & 0 & \cdots & 0 \\ 0 & D_{n_2}(\eta) & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & D_{n_k}(\eta) \end{bmatrix}, \text{ where } D_m(\eta) = \begin{bmatrix} \eta & 0 & \cdots & 0 \\ 0 & \eta^2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & \eta^m \end{bmatrix}.$$

Then the norm defined by

$$\| \| M \| \| \equiv \| D(\frac{1}{\varepsilon}) S^{-1} M S D(\varepsilon) \|_1$$

has the property that $\| \| A \| \| \leq \rho(A) + \varepsilon$. Define a norm on \mathbb{R}^n by $\| \| \mathbf{x} \| \|_{\mathbb{R}^n} = \| D(\frac{1}{\varepsilon}) S^{-1} \mathbf{x} \|_1$. Then

$$\begin{aligned} \| \| M \mathbf{x} \| \|_{\mathbb{R}^n} &= \| D(\frac{1}{\varepsilon}) S^{-1} M \mathbf{x} \|_1 = \| D(\frac{1}{\varepsilon}) S^{-1} M S D(\varepsilon) D(\frac{1}{\varepsilon}) S^{-1} \mathbf{x} \|_1 \\ &\leq \| D(\frac{1}{\varepsilon}) S^{-1} M S D(\varepsilon) \|_1 \| D(\frac{1}{\varepsilon}) S^{-1} \mathbf{x} \|_1 = \| \| M \| \| \| \mathbf{x} \| \|_{\mathbb{R}^n} \end{aligned}$$

which implies that $\| \| \cdot \| \|$ is an subordinate norm. □

§5.3 Norms on Vectors and Matrices

Proof (cont'd).

For each $m \in \mathbb{N}$ and $\eta > 0$, define

$$D(\eta) = \begin{bmatrix} D_{n_1}(\eta) & 0 & \cdots & 0 \\ 0 & D_{n_2}(\eta) & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & D_{n_k}(\eta) \end{bmatrix}, \text{ where } D_m(\eta) = \begin{bmatrix} \eta & 0 & \cdots & 0 \\ 0 & \eta^2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & \eta^m \end{bmatrix}.$$

Then the norm defined by

$$\| \| M \| \| \equiv \| D(\frac{1}{\varepsilon}) S^{-1} M S D(\varepsilon) \|_1$$

has the property that $\| \| A \| \| \leq \rho(A) + \varepsilon$. Define a norm on \mathbb{R}^n by $\| \mathbf{x} \|_{\mathbb{R}^n} = \| D(\frac{1}{\varepsilon}) S^{-1} \mathbf{x} \|_1$. Then

$$\begin{aligned} \| M \mathbf{x} \|_{\mathbb{R}^n} &= \| D(\frac{1}{\varepsilon}) S^{-1} M \mathbf{x} \|_1 = \| D(\frac{1}{\varepsilon}) S^{-1} M S D(\varepsilon) D(\frac{1}{\varepsilon}) S^{-1} \mathbf{x} \|_1 \\ &\leq \| D(\frac{1}{\varepsilon}) S^{-1} M S D(\varepsilon) \|_1 \| D(\frac{1}{\varepsilon}) S^{-1} \mathbf{x} \|_1 = \| \| M \| \| \| \mathbf{x} \|_{\mathbb{R}^n} \end{aligned}$$

which implies that $\| \| \cdot \| \|$ is an subordinate norm. □

§5.3 Norms on Vectors and Matrices

Definition

A square matrix A is said to be convergent (to zero matrix) if for all $1 \leq i, j \leq n$ the (i, j) -entry of A^n converges to 0 as $n \rightarrow \infty$.

Example

$$A = \begin{bmatrix} \frac{1}{2} & 0 \\ \frac{1}{4} & \frac{1}{2} \end{bmatrix} \Rightarrow A^2 = \begin{bmatrix} \frac{1}{4} & 0 \\ \frac{1}{4} & \frac{1}{4} \end{bmatrix} \Rightarrow A^3 = \begin{bmatrix} \frac{1}{8} & 0 \\ \frac{3}{16} & \frac{1}{8} \end{bmatrix} \Rightarrow \dots$$

By induction, one can show that

$$A^k = \begin{bmatrix} \left(\frac{1}{2}\right)^k & 0 \\ \frac{k}{2^{k+1}} & \left(\frac{1}{2}\right)^k \end{bmatrix}.$$

Since $\lim_{k \rightarrow \infty} \left(\frac{1}{2}\right)^k = 0$ and $\lim_{k \rightarrow \infty} \frac{k}{2^{k+1}} = 0$, A is a convergent matrix.

§5.3 Norms on Vectors and Matrices

Theorem

The following statements are equivalent:

- ① A is a convergent matrix;
- ② $\lim_{n \rightarrow \infty} \|A^n\| = 0$ for some matrix norm;
- ③ $\lim_{n \rightarrow \infty} \|A^n\| = 0$ for all matrix norms;
- ④ $\rho(A) < 1$;
- ⑤ $\lim_{n \rightarrow \infty} A^n \mathbf{x} = \mathbf{0}$ for all \mathbf{x} .

Remark:

- ② \Leftrightarrow ③ because all norms on a finite dimensional real vector space are equivalent.
- ① \Leftrightarrow ④ \Leftrightarrow ⑤ by writing A into Jordan canonical form.
- ① \Leftrightarrow ② by using the Frobenius norm.

§5.3 Norms on Vectors and Matrices

Lemma

Let A be a square matrix. If $\rho(A) < 1$, then $(I - A)^{-1}$ exists and

$$(I - A)^{-1} = I + A + A^2 + \cdots \left(:= \sum_{n=0}^{\infty} A^n \right).$$

Proof.

Since $\rho(A) < 1$, 1 is not an eigenvalue of A ; thus $(I - A)\mathbf{x} = \mathbf{0}$ has only trivial solution. Moreover, if A is $m \times m$, then for all $\mathbf{x} \in \mathbb{R}^m$,

$$\begin{aligned} (I - A) \sum_{n=0}^{\infty} A^n \mathbf{x} &= (I - A) \lim_{N \rightarrow \infty} \sum_{n=0}^N A^n \mathbf{x} = \lim_{N \rightarrow \infty} (I - A) \sum_{n=0}^N A^n \mathbf{x} \\ &= \lim_{N \rightarrow \infty} \left(\sum_{n=0}^N A^n \mathbf{x} - \sum_{n=0}^N A^{n+1} \mathbf{x} \right) \\ &= \lim_{N \rightarrow \infty} (\mathbf{x} - A^{N+1} \mathbf{x}) = \mathbf{x}; \end{aligned}$$

thus $(I - A)^{-1} \mathbf{x} = \sum_{n=0}^{\infty} A^n \mathbf{x}$. (Does $\sum_{n=0}^{\infty} A^n \mathbf{x}$ converges for all \mathbf{x} ?) \square

§5.3 Norms on Vectors and Matrices

Lemma

Let A be a square matrix. If $\rho(A) < 1$, then $(I - A)^{-1}$ exists and

$$(I - A)^{-1} = I + A + A^2 + \cdots \left(:= \sum_{n=0}^{\infty} A^n \right).$$

Proof.

Since $\rho(A) < 1$, 1 is not an eigenvalue of A ; thus $(I - A)\mathbf{x} = \mathbf{0}$ has only trivial solution. Moreover, if A is $m \times m$, then for all $\mathbf{x} \in \mathbb{R}^m$,

$$\begin{aligned} (I - A) \sum_{n=0}^{\infty} A^n \mathbf{x} &= (I - A) \lim_{N \rightarrow \infty} \sum_{n=0}^N A^n \mathbf{x} = \lim_{N \rightarrow \infty} (I - A) \sum_{n=0}^N A^n \mathbf{x} \\ &= \lim_{N \rightarrow \infty} \left(\sum_{n=0}^N A^n \mathbf{x} - \sum_{n=0}^N A^{n+1} \mathbf{x} \right) \\ &= \lim_{N \rightarrow \infty} (\mathbf{x} - A^{N+1} \mathbf{x}) = \mathbf{x}; \end{aligned}$$

thus $(I - A)^{-1} \mathbf{x} = \sum_{n=0}^{\infty} A^n \mathbf{x}$. (Does $\sum_{n=0}^{\infty} A^n \mathbf{x}$ converges for all \mathbf{x} ?) \square

§5.3 Norms on Vectors and Matrices

Lemma

Let A be a square matrix. If $\rho(A) < 1$, then $(I - A)^{-1}$ exists and

$$(I - A)^{-1} = I + A + A^2 + \cdots \left(:= \sum_{n=0}^{\infty} A^n \right).$$

Proof.

Since $\rho(A) < 1$, 1 is not an eigenvalue of A ; thus $(I - A)\mathbf{x} = \mathbf{0}$ has only trivial solution. Moreover, if A is $m \times m$, then for all $\mathbf{x} \in \mathbb{R}^m$,

$$\begin{aligned} (I - A) \sum_{n=0}^{\infty} A^n \mathbf{x} &= (I - A) \lim_{N \rightarrow \infty} \sum_{n=0}^N A^n \mathbf{x} = \lim_{N \rightarrow \infty} (I - A) \sum_{n=0}^N A^n \mathbf{x} \\ &= \lim_{N \rightarrow \infty} \left(\sum_{n=0}^N A^n \mathbf{x} - \sum_{n=0}^N A^{n+1} \mathbf{x} \right) \\ &= \lim_{N \rightarrow \infty} (\mathbf{x} - A^{N+1} \mathbf{x}) = \mathbf{x}; \end{aligned}$$

thus $(I - A)^{-1} \mathbf{x} = \sum_{n=0}^{\infty} A^n \mathbf{x}$. (Does $\sum_{n=0}^{\infty} A^n \mathbf{x}$ converge for all \mathbf{x} ?) \square

§5.4 Iterative Methods

Recall that in Chapter 3 to solve a nonlinear equation $f(x) = 0$ we introduce iterative method

$$x^{(k+1)} = g(x^{(k)}) \quad \text{for } k \in \mathbb{N} \cup \{0\} \text{ with } x^{(0)} \text{ given,}$$

where $f(x) = 0 \Leftrightarrow x = g(x)$, and the fixed-point of g is a solution of f .

The idea of solving $A\mathbf{x} = \mathbf{b}$ using the iterative method is based on the same concept:

- 1 $A\mathbf{x} = \mathbf{b} \Leftrightarrow \mathbf{x} = T\mathbf{x} + \mathbf{c}$ for some fixed matrix T and vector \mathbf{c} .
- 2 Given $\mathbf{x}^{(0)}$, $\mathbf{x}^{(k+1)} := T\mathbf{x}^{(k)} + \mathbf{c}$ for $k = 0, 1, 2, \dots$

§5.4 Iterative Methods

Let $Ax = b$ be a linear system of n equations, where $A = [a_{ij}]_{n \times n}$ and $b \in \mathbb{R}^n$. Then A can be decomposed into a diagonal component D , a lower triangular part L and an upper triangular part U :

$$A = \begin{bmatrix} a_{11} & \cdots & \cdots & a_{1n} \\ \vdots & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ a_{n1} & \cdots & \cdots & a_{nn} \end{bmatrix}, \quad D = \begin{bmatrix} a_{11} & 0 & \cdots & 0 \\ 0 & a_{22} & \ddots & 0 \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & a_{nn} \end{bmatrix}.$$

$$L = \begin{bmatrix} 0 & 0 & \cdots & 0 \\ a_{21} & 0 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ a_{n1} & \cdots & a_{n(n-1)} & 0 \end{bmatrix}, \quad U = \begin{bmatrix} 0 & a_{12} & \cdots & a_{1n} \\ 0 & 0 & \ddots & \vdots \\ \vdots & \ddots & \ddots & a_{(n-1)n} \\ 0 & \cdots & 0 & 0 \end{bmatrix}$$

- 1 Jacobi method: $Ax = b \Leftrightarrow Dx = -(L+U)x + b$.
- 2 Gauss-Seidel method: $Ax = b \Leftrightarrow (D+L)x = -Ux + b$.

§5.4 Iterative Methods

Let $A\mathbf{x} = \mathbf{b}$ be a linear system of n equations, where $A = [a_{ij}]_{n \times n}$ and $\mathbf{b} \in \mathbb{R}^n$. Then A can be decomposed into a diagonal component D , a lower triangular part L and an upper triangular part U :

$$A = \begin{bmatrix} a_{11} & \cdots & \cdots & a_{1n} \\ \vdots & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ a_{n1} & \cdots & \cdots & a_{nn} \end{bmatrix}, \quad D = \begin{bmatrix} a_{11} & 0 & \cdots & 0 \\ 0 & a_{22} & \ddots & 0 \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & a_{nn} \end{bmatrix}.$$

$$L = \begin{bmatrix} 0 & 0 & \cdots & 0 \\ a_{21} & 0 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ a_{n1} & \cdots & a_{n(n-1)} & 0 \end{bmatrix}, \quad U = \begin{bmatrix} 0 & a_{12} & \cdots & a_{1n} \\ 0 & 0 & \ddots & \vdots \\ \vdots & \ddots & \ddots & a_{(n-1)n} \\ 0 & \cdots & 0 & 0 \end{bmatrix}$$

- 1 Jacobi method: $A\mathbf{x} = \mathbf{b} \Leftrightarrow D\mathbf{x} = -(L+U)\mathbf{x} + \mathbf{b}$.
- 2 Gauss-Seidel method: $A\mathbf{x} = \mathbf{b} \Leftrightarrow (D+L)\mathbf{x} = -U\mathbf{x} + \mathbf{b}$.

§5.4 Iterative Methods

- ① The Jacobi method of solving $A\mathbf{x} = \mathbf{b}$ is the iterative method

$$\mathbf{x}^{(k+1)} = D^{-1}[\mathbf{b} - (L+U)\mathbf{x}^{(k)}] = -D^{-1}(L+U)\mathbf{x}^{(k)} + D^{-1}\mathbf{b},$$

and the element-based formula is thus

$$x_i^{(k+1)} = \frac{-\sum_{j=1, j \neq i}^n a_{ij}x_j^{(k)} + b_i}{a_{ii}} \quad \forall k \in \mathbb{N} \cup \{0\}.$$

- ② The Gauss-Seidel method of solving $A\mathbf{x} = \mathbf{b}$ is the iterative method

$$\mathbf{x}^{(k+1)} = (D+L)^{-1}[\mathbf{b} - U\mathbf{x}^{(k)}] = -(D+L)^{-1}U\mathbf{x}^{(k)} + (D+L)^{-1}\mathbf{b},$$

and the element-based formula is thus

$$x_i^{(k+1)} = \frac{-\sum_{j=1}^{i-1} a_{ij}x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij}x_j^{(k)} + b_i}{a_{ii}} \quad \forall k \in \mathbb{N} \cup \{0\}.$$

§5.4 Iterative Methods

- ① The Jacobi method of solving $A\mathbf{x} = \mathbf{b}$ is the iterative method

$$\mathbf{x}^{(k+1)} = D^{-1}[\mathbf{b} - (L+U)\mathbf{x}^{(k)}] = -D^{-1}(L+U)\mathbf{x}^{(k)} + D^{-1}\mathbf{b},$$

and the element-based formula is thus

$$x_i^{(k+1)} = \frac{-\sum_{j=1, j \neq i}^n a_{ij}x_j^{(k)} + b_i}{a_{ii}} \quad \forall k \in \mathbb{N} \cup \{0\}.$$

- ② The Gauss-Seidel method of solving $A\mathbf{x} = \mathbf{b}$ is the iterative method

$$\mathbf{x}^{(k+1)} = (D+L)^{-1}[\mathbf{b} - U\mathbf{x}^{(k)}] = -(D+L)^{-1}U\mathbf{x}^{(k)} + (D+L)^{-1}\mathbf{b},$$

and the element-based formula is thus

$$x_i^{(k+1)} = \frac{-\sum_{j=1}^{i-1} a_{ij}x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij}x_j^{(k)} + b_i}{a_{ii}} \quad \forall k \in \mathbb{N} \cup \{0\}.$$

§5.4 Iterative Methods

Example (Solving $A\mathbf{x} = \mathbf{b}$ using Jacobi and Gauss-Seidel methods)

Consider a linear system:

$$\begin{cases} 10x_1 - 1x_2 + 2x_3 + 0x_4 = 6 \\ -x_1 + 11x_2 - 1x_3 + 3x_4 = 25 \\ 2x_1 - 1x_2 + 10x_3 - 1x_4 = -11 \\ 0x_1 + 3x_2 - 1x_3 + 8x_4 = 15 \end{cases}$$

or equivalently,

$$\begin{bmatrix} 10 & -1 & 2 & 0 \\ -1 & 11 & -1 & 3 \\ 2 & -1 & 10 & -1 \\ 0 & 3 & -1 & 8 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 6 \\ 25 \\ -11 \\ 15 \end{bmatrix}.$$

Exact unique solution: $\mathbf{x} = (1, 2, -1, 1)^\top$.

§5.4 Iterative Methods

Example (cont'd)

We first rewrite the linear system as

$$x_1 = 0 + \frac{1}{10}x_2 - \frac{2}{10}x_3 + 0 + \frac{6}{10}$$

$$x_2 = \frac{1}{11}x_1 + 0 + \frac{1}{11}x_3 - \frac{3}{11}x_4 + \frac{25}{11}$$

$$x_3 = -\frac{2}{10}x_1 + \frac{1}{10}x_2 + 0 + \frac{1}{10}x_4 - \frac{11}{10}$$

$$x_4 = 0 - \frac{3}{8}x_2 + \frac{1}{8}x_3 + 0 + \frac{15}{8}$$

which, written in matrix form, is

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = T\mathbf{x} + \mathbf{c} \equiv \begin{bmatrix} 0 & \frac{1}{10} & -\frac{2}{10} & 0 \\ \frac{1}{11} & 0 & \frac{1}{11} & -\frac{3}{11} \\ -\frac{2}{10} & \frac{1}{10} & 0 & \frac{1}{10} \\ 0 & -\frac{3}{8} & \frac{1}{8} & 0 \end{bmatrix} \mathbf{x} + \begin{bmatrix} \frac{6}{10} \\ \frac{25}{11} \\ -\frac{11}{10} \\ \frac{15}{8} \end{bmatrix}.$$

§5.4 Iterative Methods

Example (cont'd)

If $\mathbf{x}^{(0)} = (0, 0, 0, 0)^\top$, then the Jacobi method provides

$$\mathbf{x}^{(1)} = T\mathbf{x}^{(0)} + \mathbf{c} = \begin{bmatrix} \frac{6}{10} \\ \frac{25}{11} \\ -\frac{11}{10} \\ \frac{15}{8} \end{bmatrix} = \begin{bmatrix} 0.6000 \\ 2.2727 \\ -1.1000 \\ 1.8750 \end{bmatrix}.$$

$$\Rightarrow \mathbf{x}^{(2)} = T\mathbf{x}^{(1)} + \mathbf{c} \Rightarrow \dots$$

$$\Rightarrow \frac{\|\mathbf{x}^{(10)} - \mathbf{x}^{(9)}\|_\infty}{\|\mathbf{x}^{(10)}\|_\infty} \approx \frac{8.0 \times 10^{-4}}{1.9998} < 10^{-3} \quad \text{stop! (Stopping criteria)}$$

$$\Rightarrow \mathbf{x} \approx \mathbf{x}^{(10)} \approx \begin{bmatrix} 1.00011860 \\ 1.99976795 \\ -0.99982814 \\ 0.99978598 \end{bmatrix}.$$

§5.4 Iterative Methods

Example (cont'd)

For the Gauss-Seidel method, we let $\mathbf{x}^{(0)} = (0, 0, 0, 0)^\top$ and for $k = 0, 1, 2, \dots$ define

$$x_1^{(k+1)} = 0 + \frac{1}{10}x_2^{(k)} - \frac{2}{10}x_3^{(k)} + 0 + \frac{6}{10}$$

$$x_2^{(k+1)} = \frac{1}{11}x_1^{(k+1)} + 0 + \frac{1}{11}x_3^{(k)} - \frac{3}{11}x_4^{(k)} + \frac{25}{11}$$

$$x_3^{(k+1)} = -\frac{2}{10}x_1^{(k+1)} + \frac{1}{10}x_2^{(k+1)} + 0 + \frac{1}{10}x_4^{(k)} - \frac{11}{10}$$

$$x_4^{(k+1)} = 0 - \frac{3}{8}x_2^{(k+1)} + \frac{1}{8}x_3^{(k+1)} + 0 + \frac{15}{8}$$

- Need to proceed from the top line to the bottom line:

Solving for $x_1^{(k+1)}$ from the first equation, and then use this solution to solve $x_2^{(k+1)}$ from the second equation, and so on.

$$\Rightarrow \frac{\|\mathbf{x}^{(5)} - \mathbf{x}^{(4)}\|_\infty}{\|\mathbf{x}^{(5)}\|_\infty} = 4.0 \times 10^{-4} < 10^{-3} \quad \text{stop!} \quad \mathbf{x} \approx \mathbf{x}^{(5)}.$$

§5.4 Iterative Methods

Theorem

Let T be an $n \times n$ real matrix. For any $\mathbf{x}^{(0)} \in \mathbb{R}^n$, the sequence $\{\mathbf{x}^{(k)}\}_{k=1}^{\infty}$ defined by

$$\mathbf{x}^{(k+1)} := T\mathbf{x}^{(k)} + \mathbf{c}, \quad k \in \mathbb{N} \cup \{0\},$$

converges to the unique solution of $\mathbf{x} = T\mathbf{x} + \mathbf{c}$ if and only if $\rho(T) < 1$.

Proof.

(\Leftarrow) Since $\rho(T) < 1$, $(I - T)^{-1}$ exists; thus $\mathbf{x} = T\mathbf{x} + \mathbf{c}$ has a unique solution. Moreover, there exists a subordinate matrix norm $\|\cdot\|$ and a norm $\|\cdot\|$ on \mathbb{R}^n such that $\|T\| < 1$ and $\|T\mathbf{x}\| \leq \|T\|\|\mathbf{x}\|$ for all $\mathbf{x} \in \mathbb{R}^n$. Therefore, the mapping $\mathbf{x} \mapsto T\mathbf{x} + \mathbf{c}$ is a contraction mapping, and the contraction mapping principle implies that the sequence $\{\mathbf{x}^{(k)}\}_{k=1}^{\infty}$ defined by $\mathbf{x}^{(k+1)} = T\mathbf{x}^{(k)} + \mathbf{c}$ converges (to the solution of $\mathbf{x} = T\mathbf{x} + \mathbf{c}$). \square

§5.4 Iterative Methods

Theorem

Let T be an $n \times n$ real matrix. For any $\mathbf{x}^{(0)} \in \mathbb{R}^n$, the sequence $\{\mathbf{x}^{(k)}\}_{k=1}^{\infty}$ defined by

$$\mathbf{x}^{(k+1)} := T\mathbf{x}^{(k)} + \mathbf{c}, \quad k \in \mathbb{N} \cup \{0\},$$

converges to the unique solution of $\mathbf{x} = T\mathbf{x} + \mathbf{c}$ if and only if $\rho(T) < 1$.

Proof.

(\Leftarrow) Since $\rho(T) < 1$, $(I - T)^{-1}$ exists; thus $\mathbf{x} = T\mathbf{x} + \mathbf{c}$ has a unique solution. Moreover, there exists a subordinate matrix norm $\| \cdot \|$ and a norm $\| \cdot \|$ on \mathbb{R}^n such that $\|T\| < 1$ and $\|T\mathbf{x}\| \leq \|T\| \|\mathbf{x}\|$ for all $\mathbf{x} \in \mathbb{R}^n$. Therefore, the mapping $\mathbf{x} \mapsto T\mathbf{x} + \mathbf{c}$ is a contraction mapping, and the contraction mapping principle implies that the sequence $\{\mathbf{x}^{(k)}\}_{k=1}^{\infty}$ defined by $\mathbf{x}^{(k+1)} = T\mathbf{x}^{(k)} + \mathbf{c}$ converges (to the solution of $\mathbf{x} = T\mathbf{x} + \mathbf{c}$). \square

§5.4 Iterative Methods

Theorem

Let T be an $n \times n$ real matrix. For any $\mathbf{x}^{(0)} \in \mathbb{R}^n$, the sequence $\{\mathbf{x}^{(k)}\}_{k=1}^{\infty}$ defined by

$$\mathbf{x}^{(k+1)} := T\mathbf{x}^{(k)} + \mathbf{c}, \quad k \in \mathbb{N} \cup \{0\},$$

converges to the unique solution of $\mathbf{x} = T\mathbf{x} + \mathbf{c}$ if and only if $\rho(T) < 1$.

Proof.

(\Leftarrow) Since $\rho(T) < 1$, $(I - T)^{-1}$ exists; thus $\mathbf{x} = T\mathbf{x} + \mathbf{c}$ has a unique solution. Moreover, there exists a subordinate matrix norm $\| \cdot \|$ and a norm $\| \cdot \|$ on \mathbb{R}^n such that $\|T\| < 1$ and $\|T\mathbf{x}\| \leq \|T\| \|\mathbf{x}\|$ for all $\mathbf{x} \in \mathbb{R}^n$. Therefore, the mapping $\mathbf{x} \mapsto T\mathbf{x} + \mathbf{c}$ is a contraction mapping, and the contraction mapping principle implies that the sequence $\{\mathbf{x}^{(k)}\}_{k=1}^{\infty}$ defined by $\mathbf{x}^{(k+1)} = T\mathbf{x}^{(k)} + \mathbf{c}$ converges (to the solution of $\mathbf{x} = T\mathbf{x} + \mathbf{c}$). \square

§5.4 Iterative Methods

Theorem

Let T be an $n \times n$ real matrix. For any $\mathbf{x}^{(0)} \in \mathbb{R}^n$, the sequence $\{\mathbf{x}^{(k)}\}_{k=1}^{\infty}$ defined by

$$\mathbf{x}^{(k+1)} := T\mathbf{x}^{(k)} + \mathbf{c}, \quad k \in \mathbb{N} \cup \{0\},$$

converges to the unique solution of $\mathbf{x} = T\mathbf{x} + \mathbf{c}$ if and only if $\rho(T) < 1$.

Proof.

(\Rightarrow) Let $\mathbf{z} \in \mathbb{R}^n$ be given, and \mathbf{x} be the unique solution to $\mathbf{x} = T\mathbf{x} + \mathbf{c}$. Define $\mathbf{x}^{(0)} = \mathbf{x} - \mathbf{z}$. Then

$$\mathbf{x}^{(1)} = T\mathbf{x}^{(0)} + \mathbf{c} = T\mathbf{x} - T\mathbf{z} + \mathbf{c} = \mathbf{x} - T\mathbf{z}$$

which further implies

$$\mathbf{x}^{(2)} = T\mathbf{x}^{(1)} + \mathbf{c} = T\mathbf{x} - T^2\mathbf{z} + \mathbf{c} = \mathbf{x} - T^2\mathbf{z}.$$

By induction, $\mathbf{x}^{(k)} = \mathbf{x} - T^k\mathbf{z}$. Since $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}$, we have

$\lim_{k \rightarrow \infty} T^k\mathbf{z} = \mathbf{0}$. Then $\rho(T) < 1$ due to the previous theorem. \square

§5.4 Iterative Methods

Corollary

- ① Let $\mathbf{x}^{(0)} \in \mathbb{R}^n$, and $\{\mathbf{x}^{(k)}\}_{k=1}^{\infty}$ be a sequence defined by $\mathbf{x}^{(k+1)} := T\mathbf{x}^{(k)} + \mathbf{c}$, $k \geq 0$. If $\|T\| < 1$ for some natural matrix norm, then $\{\mathbf{x}^{(k)}\}_{k=1}^{\infty}$ converges to the unique solution of $\mathbf{x} = T\mathbf{x} + \mathbf{c}$ and

- $\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \|T\|^k \|\mathbf{x} - \mathbf{x}^{(0)}\|.$
- $\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \frac{\|T\|^k}{1 - \|T\|} \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\|.$

- ② If A is *strictly diagonally dominant*, then for any $\mathbf{x}^{(0)} \in \mathbb{R}^n$, both the Jacobi and Gauss-Seidel methods give sequences $\{\mathbf{x}^{(k)}\}_{k=1}^{\infty}$ that converge to the unique solution of $A\mathbf{x} = \mathbf{b}$.

§5.4 Iterative Methods

Successive Over Relaxation (SOR):

① The Gauss-Seidel method:

$$x_i^{(k)} = \frac{1}{a_{ii}} \left[- \sum_{j=1}^{i-1} a_{ij} x_j^{(k)} - \sum_{j=i+1}^n a_{ij} x_j^{(k-1)} + b_i \right].$$

② Successive over-relaxation: for $\omega > 0$,

$$x_i^{(k)} = (1-\omega)x_i^{(k-1)} + \frac{\omega}{a_{ii}} \left[- \sum_{j=1}^{i-1} a_{ij} x_j^{(k)} - \sum_{j=i+1}^n a_{ij} x_j^{(k-1)} + b_i \right].$$

$$\Leftrightarrow a_{ii} x_i^{(k)} + \omega \sum_{j=1}^{i-1} a_{ij} x_j^{(k)} = (1-\omega) a_{ii} x_i^{(k-1)} - \omega \sum_{j=i+1}^n a_{ij} x_j^{(k-1)} + \omega b_i$$

$$\Leftrightarrow (D + \omega L) \mathbf{x}^{(k)} = \left[(1-\omega) D - \omega U \right] \mathbf{x}^{(k-1)} + \omega \mathbf{b}$$

$$\Leftrightarrow \mathbf{x}^{(k)} = (D + \omega L)^{-1} \left[(1-\omega) D - \omega U \right] \mathbf{x}^{(k-1)} + \omega (D + \omega L)^{-1} \mathbf{b}$$

$$\Leftrightarrow \mathbf{x}^{(k)} = T_\omega \mathbf{x}^{(k-1)} + \mathbf{c}_\omega$$

§5.4 Iterative Methods

Gauss-Seidel:

$$\mathbf{x}^{(k+1)} = -(D+L)^{-1}U\mathbf{x}^{(k)} + (D+L)^{-1}\mathbf{b}.$$

SOR:

$$\mathbf{x}^{(k)} = (D + \omega L)^{-1} \left[(1 - \omega)D - \omega U \right] \mathbf{x}^{(k-1)} + \omega(D + \omega L)^{-1} \mathbf{b}.$$

Different parameter ω can be chosen according to the need. In general,

- $\omega = 1$: the Gauss-Seidel method.
- $0 < \omega < 1$: when Gauss-Seidel diverges.
- $\omega > 1$: when Gauss-Seidel converges.

§5.4 Iterative Methods

Example

Consider a linear system

$$\begin{cases} 4x_1 + 3x_2 + 0 & = 24 \\ 3x_1 + 4x_2 - x_3 & = 30 \\ 0 - x_2 + 4x_3 & = -24 \end{cases}$$

Exact unique solution: $x = (3, 4, -5)^\top$.

- ① Let $x^{(0)} = (1, 1, 1)^\top$. The Gauss-Seidel method:

$$\begin{cases} x_1^{(k)} = -0.75x_2^{(k-1)} + 6 \\ x_2^{(k)} = -0.75x_1^{(k)} + 0.25x_3^{(k-1)} + 7.5 \\ x_3^{(k)} = 0.25x_2^{(k)} - 6 \end{cases}$$

- ② Let $x^{(0)} = (1, 1, 1)^\top$. The SOR with $\omega = 1.25$:

$$\begin{cases} x_1^{(k)} = -0.25x_1^{(k-1)} - 0.9375x_2^{(k-1)} + 7.5 \\ x_2^{(k)} = -0.9375x_1^{(k)} - 0.25x_2^{(k-1)} + 0.3125x_3^{(k-1)} + 9.375 \\ x_3^{(k)} = 0.3125x_2^{(k)} - 0.25x_3^{(k-1)} - 7.5 \end{cases}$$

§5.4 Iterative Methods

Theorem

- 1 If $a_{ii} \neq 0$ for all $i = 1, 2, \dots, n$, then $\rho(T_\omega) \geq |\omega - 1|$. This implies the SOR method can converge **only if** $0 < \omega < 2$.
- 2 If A is symmetric positive definite and $0 < \omega < 2$, then the SOR method converges for any $\mathbf{x}^{(0)}$.

§5.5 Absolute Error, Relative Error and Condition Number

- ① Suppose that we want to solve the linear system $A\mathbf{x} = \mathbf{b}$, but \mathbf{b} is somehow perturbed to $\tilde{\mathbf{b}}$ (this may happen when we convert a real \mathbf{b} to a floating-point \mathbf{b}).

- ② Then actual solution would satisfy a slightly different linear system

$$A\tilde{\mathbf{x}} = \tilde{\mathbf{b}}.$$

- ③ **Question:** Is $\tilde{\mathbf{x}}$ very different from the desired solution \mathbf{x} of the original system?
- ④ Of course, the answer should depend on **how good the matrix A is.**
- ⑤ Let $\|\cdot\|$ be a vector norm, we consider two types of errors:
- **absolute error:** $\|\mathbf{x} - \tilde{\mathbf{x}}\|$
 - **relative error:** $\|\mathbf{x} - \tilde{\mathbf{x}}\| / \|\mathbf{x}\|$

§5.5 Absolute Error, Relative Error and Condition Number

- For the absolute error, we have

$$\|\mathbf{x} - \tilde{\mathbf{x}}\| = \|A^{-1}\mathbf{b} - A^{-1}\tilde{\mathbf{b}}\| = \|A^{-1}(\mathbf{b} - \tilde{\mathbf{b}})\| \leq \|A^{-1}\| \|\mathbf{b} - \tilde{\mathbf{b}}\|.$$

Therefore, the absolute error of \mathbf{x} depends on two factors: **the absolute error of \mathbf{b}** and **the matrix norm of A^{-1}** .

- For the relative error, we have

$$\begin{aligned} \|\mathbf{x} - \tilde{\mathbf{x}}\| &= \|A^{-1}\mathbf{b} - A^{-1}\tilde{\mathbf{b}}\| = \|A^{-1}(\mathbf{b} - \tilde{\mathbf{b}})\| \\ &\leq \|A^{-1}\| \|\mathbf{b} - \tilde{\mathbf{b}}\| = \|A^{-1}\| \|A\mathbf{x}\| \frac{\|\mathbf{b} - \tilde{\mathbf{b}}\|}{\|\mathbf{b}\|} \\ &\leq \|A^{-1}\| \|A\| \|\mathbf{x}\| \frac{\|\mathbf{b} - \tilde{\mathbf{b}}\|}{\|\mathbf{b}\|}; \end{aligned}$$

that is

$$\frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|}{\|\mathbf{x}\|} \leq \|A^{-1}\| \|A\| \frac{\|\mathbf{b} - \tilde{\mathbf{b}}\|}{\|\mathbf{b}\|}.$$

Therefore, the relative error of \mathbf{x} depends on two factors: **the relative error of \mathbf{b}** and **$\|A\| \|A^{-1}\|$** .

§5.5 Absolute Error, Relative Error and Condition Number

- For the absolute error, we have

$$\|\mathbf{x} - \tilde{\mathbf{x}}\| = \|A^{-1}\mathbf{b} - A^{-1}\tilde{\mathbf{b}}\| = \|A^{-1}(\mathbf{b} - \tilde{\mathbf{b}})\| \leq \|A^{-1}\| \|\mathbf{b} - \tilde{\mathbf{b}}\|.$$

Therefore, the absolute error of \mathbf{x} depends on two factors: **the absolute error of \mathbf{b}** and **the matrix norm of A^{-1}** .

- For the relative error, we have

$$\begin{aligned} \|\mathbf{x} - \tilde{\mathbf{x}}\| &= \|A^{-1}\mathbf{b} - A^{-1}\tilde{\mathbf{b}}\| = \|A^{-1}(\mathbf{b} - \tilde{\mathbf{b}})\| \\ &\leq \|A^{-1}\| \|\mathbf{b} - \tilde{\mathbf{b}}\| = \|A^{-1}\| \|A\mathbf{x}\| \frac{\|\mathbf{b} - \tilde{\mathbf{b}}\|}{\|\mathbf{b}\|} \\ &\leq \|A^{-1}\| \|A\| \|\mathbf{x}\| \frac{\|\mathbf{b} - \tilde{\mathbf{b}}\|}{\|\mathbf{b}\|}; \end{aligned}$$

that is

$$\frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|}{\|\mathbf{x}\|} \leq \|A^{-1}\| \|A\| \frac{\|\mathbf{b} - \tilde{\mathbf{b}}\|}{\|\mathbf{b}\|}.$$

Therefore, the relative error of \mathbf{x} depends on two factors: **the relative error of \mathbf{b}** and **$\|A\| \|A^{-1}\|$** .

§5.5 Absolute Error, Relative Error and Condition Number

Definition

For a given subordinate matrix norm $\|\cdot\|$, the condition number of the matrix A is the number

$$\kappa(A) := \|A\| \|A^{-1}\|.$$

$\kappa(A)$ measures how good the matrix A is.

Example

Let $\varepsilon > 0$ and

$$A = \begin{bmatrix} 1 & 1 + \varepsilon \\ 1 - \varepsilon & 1 \end{bmatrix} \Rightarrow A^{-1} = \varepsilon^{-2} \begin{bmatrix} 1 & -1 - \varepsilon \\ -1 + \varepsilon & 1 \end{bmatrix}.$$

Then $\|A\|_{\infty} = 2 + \varepsilon$, $\|A^{-1}\|_{\infty} = \varepsilon^{-2}(2 + \varepsilon)$, and

$$\kappa(A) = \left(\frac{2 + \varepsilon}{\varepsilon}\right)^2 \geq \frac{4}{\varepsilon^2}.$$

§5.5 Absolute Error, Relative Error and Condition Number

① For example, if $\varepsilon = 0.01$, then $\kappa(A) \geq 40000$.

② What does this mean?

It means that the relative error in \mathbf{x} can be 40000 times greater than the relative error in b .

③ If $\kappa(A)$ is large, we say that A is **ill-conditioned**, otherwise A is **well-conditioned**.

④ In the ill-conditioned case, the solution is very sensitive to the small changes in the right-hand vector b (higher precision in b may be needed).

§5.5 Absolute Error, Relative Error and Condition Number

Consider the linear system $A\mathbf{x} = \mathbf{b}$. Let $\tilde{\mathbf{x}}$ be a computed solution (which is an approximation to \mathbf{x}). We define

- ① Residual vector: $\mathbf{r} = \mathbf{b} - A\tilde{\mathbf{x}}$.
- ② Error vector: $\mathbf{e} = \mathbf{x} - \tilde{\mathbf{x}}$.

Then $A\mathbf{e} = A\mathbf{x} - A\tilde{\mathbf{x}} = \mathbf{b} - A\tilde{\mathbf{x}} = \mathbf{r}$.

Theorem (bounds involving condition number)

Let A be a square matrix, \mathbf{x} be the solution of $A\mathbf{x} = \mathbf{b}$, and \mathbf{r} , \mathbf{e} are the residual vector and the error vector associated with a computed solution $\tilde{\mathbf{x}}$, respectively. Then

$$\frac{1}{\kappa(A)} \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} \leq \frac{\|\mathbf{e}\|}{\|\mathbf{x}\|} \leq \kappa(A) \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|}.$$

§5.5 Absolute Error, Relative Error and Condition Number

Theorem (bounds involving condition number)

Let A be a square matrix, \mathbf{x} be the solution of $A\mathbf{x} = \mathbf{b}$, and \mathbf{r} , \mathbf{e} are the residual vector and the error vector associated with a computed solution $\tilde{\mathbf{x}}$, respectively. Then

$$\frac{1}{\kappa(A)} \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} \leq \frac{\|\mathbf{e}\|}{\|\mathbf{x}\|} \leq \kappa(A) \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|}.$$

Proof.

Since $A\mathbf{e} = \mathbf{r}$, $\mathbf{e} = A^{-1}\mathbf{r}$; thus

$$\|\mathbf{e}\| \|\mathbf{b}\| = \|A^{-1}\mathbf{r}\| \|A\mathbf{x}\| \leq \|A^{-1}\| \|\mathbf{r}\| \|A\| \|\mathbf{x}\| = \kappa(A) \|\mathbf{r}\| \|\mathbf{x}\|$$

which further implies that $\frac{\|\mathbf{e}\|}{\|\mathbf{x}\|} \leq \kappa(A) \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|}$.

On the other hand, we have

$$\|\mathbf{r}\| \|\mathbf{x}\| = \|A\mathbf{e}\| \|A^{-1}\mathbf{b}\| \leq \|A\| \|\mathbf{e}\| \|A^{-1}\| \|\mathbf{b}\| = \kappa(A) \|\mathbf{e}\| \|\mathbf{b}\|$$

which shows that $\frac{1}{\kappa(A)} \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} \leq \frac{\|\mathbf{e}\|}{\|\mathbf{x}\|}$. □

§5.5 Absolute Error, Relative Error and Condition Number

Theorem (bounds involving condition number)

Let A be a square matrix, \mathbf{x} be the solution of $A\mathbf{x} = \mathbf{b}$, and \mathbf{r} , \mathbf{e} are the residual vector and the error vector associated with a computed solution $\tilde{\mathbf{x}}$, respectively. Then

$$\frac{1}{\kappa(A)} \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} \leq \frac{\|\mathbf{e}\|}{\|\mathbf{x}\|} \leq \kappa(A) \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|}.$$

Proof.

Since $A\mathbf{e} = \mathbf{r}$, $\mathbf{e} = A^{-1}\mathbf{r}$; thus

$$\|\mathbf{e}\| \|\mathbf{b}\| = \|A^{-1}\mathbf{r}\| \|A\mathbf{x}\| \leq \|A^{-1}\| \|\mathbf{r}\| \|A\| \|\mathbf{x}\| = \kappa(A) \|\mathbf{r}\| \|\mathbf{x}\|$$

which further implies that $\frac{\|\mathbf{e}\|}{\|\mathbf{x}\|} \leq \kappa(A) \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|}$.

On the other hand, we have

$$\|\mathbf{r}\| \|\mathbf{x}\| = \|A\mathbf{e}\| \|A^{-1}\mathbf{b}\| \leq \|A\| \|\mathbf{e}\| \|A^{-1}\| \|\mathbf{b}\| = \kappa(A) \|\mathbf{e}\| \|\mathbf{b}\|$$

which shows that $\frac{1}{\kappa(A)} \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} \leq \frac{\|\mathbf{e}\|}{\|\mathbf{x}\|}$. □