

Prelim Notes for Numerical Analysis

Wenqiang Feng
wfeng1@utk.edu

Department of Mathematics,
University of Tennessee, Knoxville, TN, 37996

April 21, 2019

Abstract

This note is for my Numerical Analysis prelim exam in University of Tennessee at Knoxville. This note is intended to assist my prelim examination preparation. You can download and distribute it. [Please be aware, however, that the note might contain typos as well as incorrect or inaccurate solutions](#). At here, I also would like to thank Ligu Wang for their help in some problems.

Contents

List of Figures	1
List of Tables	1
1 Preliminaries	2
1.1 Linear Algebra Preliminaries	2
1.1.1 Common Properties	2
1.1.2 Similar and diagonalization	4
1.1.3 Eigenvalues and Eigenvectors	5
1.1.4 Unitary matrices	6
1.1.5 Hermitian matrices	7
1.1.6 Positive definite matrices	8
1.1.7 Normal matrices	9

4	1.1.8	Common Theorems	9
5	1.2	Calculus Preliminaries	10
1	1.3	Norms' Preliminaries	14
2	1.3.1	Vector Norms	14
3	1.3.2	Matrix Norms	15
4	2	Direct Method	17
5	2.1	For square or rectangular matrices $A \in \mathbb{C}^{m,n}, m \geq n$	17
6	2.1.1	Singular Value Decomposition	17
7	2.1.2	QR Decomposition	17
8	2.2	For square matrices $A \in \mathbb{C}^{n,n}$	17
9	2.2.1	LU Decomposition	17
10	2.2.2	Cholesky Decomposition	17
1	3	Iterative Method	17
2	3.1	General Iterative Scheme	17
3	3.2	Jacobi Method	17
4	3.3	Gauss-Seidel Method	17
5	3.4	Richardson Method	17
6	3.5	Successive Over Relaxation (SOR) Method	17
7	3.6	Minimal Correction Method	17
8	3.7	Steepest Descent Method	17
9	3.8	Conjugate Gradients Method	17
10	4	Eigenvalue Problems	17
11	5	Euler Method	17
12	5.1	Euler's method	18
13	5.2	Trapezoidal Method	21
14	5.3	Theta Method	23
15	5.4	Midpoint Rule Method	23
16	6	Multistep Method	25
1	6.1	The Adams Method	25
2	6.2	The Order and Convergence of Multistep Methods	26
3	6.3	Method of A-stable verification for Multistep Methods	27
4	7	Runge-Kutta Methods	27
5	7.1	Quadrature Formulas	27
6	7.2	Explicit Runge-Kutta Formulas	27
7	7.3	Implicit Runge-Kutta Formulas	27
8	7.4	Collocation Runge-Kutta Formulas	27
9	7.5	Method of A-stable verification for Runge-Kutta Method	27
10	7.6	Problems	28
11		Appendices	31

A	Lecture notes	31
B	Trigonometric formula tables	31
C	Trigonometric tables	31

List of Figures

List of Tables

1 Preliminaries

1.1 Linear Algebra Preliminaries

1.1.1 Common Properties

Properties 1.1. (Structure of Matrices) Let $A = [A_{ij}]$ be a square or rectangular matrix, A is called

- *diagonal* : if $a_{ij} = 0, \forall i \neq j$,
- *upper triangular* : if $a_{ij} = 0, \forall i > j$,
- *upper Hessenberg* : if $a_{ij} = 0, \forall i > j + 1$,
- *block diagonal* : $A = \text{diag}(A_{11}, A_{22}, \dots, A_{nn})$,
- *tridiagonal* : if $a_{ij} = 0, \forall |i - j| > 1$,
- *lower triangular* : if $a_{ij} = 0, \forall i < j$,
- *lower Hessenberg* : if $a_{ij} = 0, \forall j > i + 1$,
- *block diagonal* : $A = \text{diag}(A_{i,i-1}, A_{ii}, \dots, A_{i,i+1})$.

Properties 1.2. (Type of Matrices) Let $A = [A_{ij}]$ be a square or rectangular matrix, A is called

- *Hermitian* : if $A^* = A$,
- *symmetric* : if $A^T = A$,
- *normal* : if $A^T A = A A^T$, when $A \in \mathbb{R}^{n \times n}$, if $A^* A = A A^*$, when $A \in \mathbb{C}^{n \times n}$,
- *skew hermitian* : if $A^* = -A$,
- *skew symmetric* : if $A^T = -A$,
- *orthogonal* : if $A^T A = I$, when $A \in \mathbb{R}^{n \times n}$, *unitary* : if $A^* A = I$, when $A \in \mathbb{C}^{n \times n}$.

Properties 1.3. (Properties of invertible matrices) Let A be $n \times n$ square matrix. If A is invertible, then

- $\det(A) \neq 0$,
- $\text{rank}(A) = n$,
- $Ax = b$ has a unique solution for every $b \in \mathbb{R}^n$
- the row vectors are linearly independent,
- the row vectors of A form a basis for \mathbb{R}^n .
- the row vectors of A span \mathbb{R}^n .
- $\text{nullity}(A) = 0$,
- $\lambda_i \neq 0$, (λ_i eigenvalues),
- $Ax = 0$ has only trivial solution,
- the column vectors are linearly independent,
- the column vectors of A form a basis for \mathbb{R}^n ,
- the column vectors of A span \mathbb{R}^n .

Properties 1.4. (Properties of conjugate transpose) Let A, B be $n \times n$ square matrix and γ be a complex constant, then

- $(A^*)^* = A$,
- $(AB)^* = B^*A^*$,
- $(A + B)^* = A^* + B^*$,
- $\det(A^*) = \det(A)$
- $\text{tr}(A^*) = \text{tr}(A)$
- $(\gamma A)^* = \gamma^* A^*$.

Properties 1.5. (Properties of similar matrices) If $A \sim B$, then

- $\det(A) = \det(B)$,
- $\text{eig}(A) = \text{eig}(B)$,
- $A \sim A$,
- $\text{rank}(A) = \text{rank}(B)$,
- if $B \sim C$, then $A \sim C$
- $B \sim A$

Properties 1.6. (Properties of Unitary Matrices) Let A be a $n \times n$ Unitary matrix, then

- $A^* = A^{-1}$,
- A^* is unitary,
- A is diagonalizable,
- A is unitarily similar to a diagonal matrix,
- the row vectors of A form an orthonormal set,
- $A^* = I$,
- A is an isometry.
- the column vectors of A form an orthonormal set.

Properties 1.7. (Properties of Hermitian Matrices) Let A be a $n \times n$ Hermitian matrix, then

- its eigenvalues are real ,
- $v_i^* v_j = 0, i \neq j$, v_i, v_j eigenvectors,
- A is unitarily diagonalizable (Spectral theorem),
- $A = H + K$, H is Hermitian and K is skew-Hermitian,

Properties 1.8. (Properties of positive definite Matrices) Let $A \in \mathbb{C}^{n \times n}$ be a positive definite Matrix and $B \in \mathbb{C}^{n \times n}$, then

- $\sigma(A) \subset (0, \infty)$,
- A is invertible,
- if B is invertible, $B^* B$ positive semidefinite ,
- if B is positive semidefinite then $\text{diag}(A) \geq 0$,
- if B is positive definite then $\text{diag}(A) > 0$.
- $B^* B$ is positive semidefinite

Properties 1.9. (Properties of determinants) Let A, B be $n \times n$ square matrix and α be a real constant, then

- $\det(A^T) = \det(A)$,
- $\det(AB) = \det(A)\det(B)$,
- $\det(\alpha A) = \alpha^n \det(A)$,
- $\det(A^{-1}) = \frac{1}{\det(A)} = \det(A)^{-1}$.

Properties 1.10. (Properties of inverse) Let A, B be $n \times n$ square matrix and α be a real constant, then

- $(A^*)^{-1} = (A^{-1})^*$,
- $(A^{-1})^{-1} = A$,
- $(AB)^{-1} = B^{-1}A^{-1}$,
- $(\alpha A)^{-1} = \frac{1}{\alpha}A^{-1}$
- $A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$, $A^{-1} = \frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$.

Properties 1.11. (Properties of Rank) Let A be $m \times n$ matrix, B be $n \times m$ matrix and P, Q are invertible $n \times n$ matrices, then

- $\text{rank}(A) \leq \min\{m, n\}$,
- $\text{rank}(A) = \text{rank}(A^*)$,
- $\text{rank}(A) + \dim(\ker(A)) = n$,
- $\text{rank}(AQ) = \text{Rank}(A) = \text{Rank}(PA)$,
- $\text{rank}(PAQ) = \text{Rank}(A)$,
- $\text{rank}(AB) \geq \text{rank}(A) + \text{rank}(B) - n$,
- $\text{rank}(AB) \leq \min\{\text{rank}(A), \text{rank}(B)\}$,
- $\text{rank}(AB) \leq \text{rank}(A) + \text{rank}(B)$.

1.1.2 Similar and diagonalization

Theorem 1.1. (Similar) *A is said to be similar to B, if there is a nonsingular matrix X, such that*

$$A = XBX^{-1}, (A \sim B).$$

Theorem 1.2. (Diagonalizable^a) *A matrix is diagonalizable, if and only if there exist a nonsingular matrix X and a diagonal matrix D such that $A = XDX^{-1}$.*

^aBeing diagonalizable has nothing to do with being invertible.

Theorem 1.3. (Diagonalizable) *A matrix is diagonalizable, if and only if all its eigenvalues are semisimple.*

Theorem 1.4. (Diagonalizable) *Suppose $\dim(A) = n$. A is said to be diagonalizable, if and only if A has n linearly independent eigenvectors.*

Corollary 1.1. (Sample question #2, summer, 2013) *Suppose $\dim(A) = n$. If A has n distinct eigenvalues, then A is diagonalizable.*

Proof. (Sketch) Suppose $n = 2$, and let λ_1, λ_2 be distinct eigenvalues of A with corresponding eigenvectors v_1, v_2 . Now, we will use contradiction to show v_1, v_2 are linearly independent. Suppose v_1, v_2 are linearly dependent, then

$$c_1 v_1 + c_2 v_2 = 0, \quad (1)$$

with c_1, c_2 are not both 0. Multiplying A on both sides of (1), then

$$c_1 A v_1 + c_2 A v_2 = c_1 \lambda_1 v_1 + c_2 \lambda_2 v_2 = 0. \quad (2)$$

Multiplying λ_1 on both sides of (1), then

$$c_1 \lambda_1 v_1 + c_2 \lambda_1 v_2 = 0. \quad (3)$$

Subtracting (3) from (2), then

$$c_2 (\lambda_2 - \lambda_1) v_2 = 0. \quad (4)$$

Since $\lambda_1 \neq \lambda_2$ and $v_2 \neq 0$, then $c_2 = 0$. Similarly, we can get $c_1 = 0$. Hence, we get the contradiction.

A similar argument gives the result for n . Then we get A has n linearly independent eigenvectors. \square

Theorem 1.5. (Diagonalizable) *Every Hermitian matrix is diagonalizable, In particular, every real symmetric matrix is diagonalizable.*

1.1.3 Eigenvalues and Eigenvectors

Theorem 1.6. if λ is an eigenvalue of A , then $\bar{\lambda}$ is an eigenvalue of A^* .

Theorem 1.7. The eigenvalues of a triangular matrix are the entries on its main diagonal.

Theorem 1.8. Let A be square matrix with eigenvalue λ and the corresponding eigenvector x .

- $\lambda^n, n \in \mathbb{Z}$ is an eigenvalue of A^n with corresponding eigenvector x ,
- if A is invertible, then $1/\lambda$ is an eigenvalue of A^{-1} with corresponding eigenvector x .

Theorem 1.9. Let A be $n \times n$ square matrix and let $\lambda_1, \lambda_2, \dots, \lambda_m$ be distinct eigenvalues of A with corresponding eigenvectors v_1, v_2, \dots, v_m . Then v_1, v_2, \dots, v_m are linear independent.

1.1.4 Unitary matrices

Definition 1.1. (Unitary Matrix) A matrix $A \in \mathbb{C}^{n \times n}$ is said to be *unitary*^a, if

$$A^* A = I.$$

^aA matrix $A \in \mathbb{R}^{n \times n}$ is said to be *orthogonal*, if

$$A^T A = I.$$

Theorem 1.10. (Angle preservation) A matrix is *unitary*, then the transformation defined by A preserves angles.

Proof. For any vectors $x, y \in \mathbb{C}^n$ that is angle θ is determined from the inner product via $\cos \theta = \frac{\langle x, y \rangle}{\|x\| \|y\|}$. Since A is unitary (and thus an isometry), then

$$\langle Ax, Ay \rangle = \langle A^* Ax, y \rangle = \langle x, y \rangle.$$

This proves the Angle preservation. □

Theorem 1.11. (Angle preservation) A matrix is real *orthogonal*, then A has the transformation form $T(\theta)$ for some θ

$$A = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} T(\theta) = \begin{bmatrix} \cos(\theta) & \sin(\theta) \\ \sin(\theta) & -\cos(\theta) \end{bmatrix} \quad (5)$$

Finally, we can easily establish the diagonalizability of the unitary matrices.

Theorem 1.12. (Shur Decomposition) A matrix $A \in \mathbb{C}^{n \times n}$ is similar to a upper triangular matrix and

$$A = UTU^{-1}, \quad (6)$$

where U is a unitary matrix, T is an upper triangular matrix.

Proof. see Appendix (A) □

Theorem 1.13. (Spectral Theorem for Unitary matrices) A is unitary, then A is diagonalizable and A is unitarily similar to a diagonal matrix.

$$A = UDU^{-1} = UDU^*, \quad (7)$$

where U is a unitary matrix, D is an diagonal matrix.

Proof. Result follows from 1.12. □

Theorem 1.14. (Spectral representation) A is unitary, then

1. A has a set of n orthogonal eigenvectors,
2. let $\{v_1, v_2, \dots, v_n\}$ be the eigenvectors w.r.t the corresponding eigenvalues $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$. The A has the representation as the sum of rank one matrices given by

$$A = \sum_{i=1}^n \lambda_i v_i v_i^T. \quad (8)$$

Note: this representation is often called the *Spectral Representation or Spectral Decomposition* of A.

Proof. see Appendix (A) □

1.1.5 Hermitian matrices

Definition 1.2. (Hermitian Matrix) A matrix is Hermitian, if

$$A^* = A.$$

Definition 1.3. Let A be Hermitian, then the spectral of A, $\sigma(A)$, is real.

Proof. Let $\lambda \in \sigma(A)$ with corresponding eigenvector v . Then

$$\langle Av, v \rangle = \langle \lambda v, v \rangle = \lambda \langle v, v \rangle \quad (9)$$

$$\langle Av, v \rangle = \langle v, A^* v \rangle = \langle v, \bar{\lambda} v \rangle = \bar{\lambda} \langle v, v \rangle. \quad (10)$$

Since $\langle v, v \rangle \neq 0$, therefore $\lambda = \bar{\lambda}$. Hence λ is real. \square

Definition 1.4. Let A be *Hermitian*, then the different eigenvector are orthogonal i.e.

$$\langle v_i, v_j \rangle = 0, i \neq j. \quad (11)$$

Proof. Let λ_1, λ_2 be the arbitrary two different eigenvalues with corresponding eigenvector v_1, v_2 . Then

$$\langle Av_1, v_2 \rangle = \langle \lambda_1 v_1, v_2 \rangle = \lambda_1 \langle v_1, v_2 \rangle \quad (12)$$

$$\langle Av_1, v_2 \rangle = \langle v_1, A^* v_2 \rangle = \langle v_1, Av_2 \rangle = \langle v_1, \lambda_2 v_2 \rangle = \lambda_2 \langle v_1, v_2 \rangle. \quad (13)$$

Since $\lambda_1 \neq \lambda_2$, therefore $\langle v_1, v_2 \rangle = 0$. \square

Theorem 1.15. (*Spectral Theorem for Hermitian matrices*) A is *Hermitian*, then A is *unitary diagonalizable*.

$$A = UDU^{-1} = UDU^*, \quad (14)$$

where U is a unitary matrix, D is an diagonal matrix.

Theorem 1.16. If A, B are *unitarily similar*, then A is *Hermitian* if and only if B is *Hermitian*.

Proof. Since A, B are *unitarily similar*, then $A = UBU^{-1}$, where U is a unitary matrix. And

$$A^* = U^{-1*} B^* U^* = U^{*-1} B^* U^* = UB^* U^{-1},$$

since U is a unitary matrix. Therefore

$$UBU^{-1} = A = A^* = UB^* U^{-1}.$$

Hence, $B = B^*$. \square

1.1.6 Positive definite matrices

Definition 1.5. (*Positive Definite Matrix*)

1. A *symmetric* real matrix $A \in \mathbb{R}^{n \times n}$ is said to be *Positive Definite*, if

$$x^T A x > 0, \quad \forall x \neq 0.$$

2. A *Hermitian* matrix $A \in \mathbb{C}^{n \times n}$ is said to be *Positive Definite*, if

$$x^* A x > 0, \quad \forall x \neq 0.$$

Theorem 1.17. Let $A, B \in \mathbb{C}^{n \times n}$. Then

1. if A is positive definite, then $\sigma(A) \subset (0, \infty)$,
2. if A is positive definite, then A is invertible,
3. B^*B is positive semidefinite,
4. if B is invertible, then B^*B is positive definite.
5. if B is positive definite, then $\text{diag}(B)$ is nonnegative,
6. if $\text{diag}(B)$ strictly positive, then B is positive definite.

Problem 1.1. (Sample question #1, summer, 2013) Suppose $A \in \mathbb{C}^{n \times n}$ is hermitian and $\sigma(A) \subset (0, \infty)$. Prove A is Hermitian Positive Defined (HPD).

Proof. Since, A is Hermitian, then is Unitary diagonalizable. i.e. $A = UDU^{-1} = UDU^*$, then

$$x^*Ax = x^*UDU^{-1}x = x^*UDU^*x = (U^*x)^*D(U^*x). \quad (15)$$

Moreover, since $\sigma(A) \subset (0, \infty)$ then $\tilde{x}^*D\tilde{x} > 0$ for any nonzero \tilde{x} . Hence

$$x^*Ax = (U^*x)^*D(U^*x) = \tilde{x}^*D\tilde{x} > 0, \text{ for any nonzero } x. \quad (16)$$

□

1.1.7 Normal matrices

Definition 1.6. (Normal Matrix) A matrix is called *normal*, if

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

Corollary 1.2. Unitary matrix and Hermitian matrix are normal matrices.

Theorem 1.18. $A \in \mathbb{C}^{n \times n}$ is normal if and only if every matrix unitarily equivalent to A is normal.

Theorem 1.19. $A \in \mathbb{C}^{n \times n}$ is normal if and only if every matrix unitarily equivalent to A is normal.

Proof. Suppose A is normal and $B = U^*AU$, where U is unitary. Then $B^*B = U^*A^*UU^*AU = U^*A^*AU = U^*AA^*U = U^*AUU^*A^*U = BB^*$, so B is normal. Conversely, If B is normal, it is easy to get that $U^*A^*AU = U^*AA^*U$, then $A^*A = AA^*$ □

Theorem 1.20. (Spectral theorem for normal matrices) If $A \in \mathbb{C}^{n \times n}$ has eigenvalues $\lambda_1, \dots, \lambda_n$, counted according to multiplicity, the following statements are equivalent.

1. A is normal,
2. A is unitarily diagonalizable,
3. $\sum_{i=1}^n \sum_{j=1}^n |a_{ij}|^2 = \sum_{i=1}^n |\lambda_i|^2$,
4. There is an orthonormal set of n eigenvectors of A .

1.1.8 Common Theorems

Definition 1.7. (Orthogonal Complement) Suppose $S \subset \mathbb{R}^n$ is a subspace. The (Orthogonal Complement) of S is defined as

$$S^\perp = \{y \in \mathbb{R}^n \mid y^T x = 0, \forall x \in S\}$$

Theorem 1.21. Suppose $A \in \mathbb{R}^{n \times n}$. Then

1. $\mathcal{R}(A)^\perp = \mathcal{N}(A^T)$,
2. $\mathcal{R}(A^T)^\perp = \mathcal{N}(A)$.

Proof. 1. For any $\tilde{y} \in \mathcal{R}(A)^\perp$, then $\tilde{y}^T y = 0, \forall y \in \mathcal{R}(A)$. And $\forall y \in \mathcal{R}(A)$, there exists x , such that $Ax = y$. Then

$$\tilde{y}^T Ax = (A^T \tilde{y})^T x = 0.$$

Since, x is arbitrary, so it must be $A^T \tilde{y} = 0$. Hence

$$\mathcal{R}(A)^\perp \subset \mathcal{N}(A^T)$$

Conversely, suppose $y \in \mathcal{N}(A^T)$, then $A^T y = 0$ and hence $(A^T y)^T x = y^T Ax = 0$ for any $x \in \mathbb{R}^n$. So, $y \in \mathcal{R}(A^T)^\perp$. Therefore

$$\mathcal{N}(A^T) \subset \mathcal{R}(A)^\perp$$

$$\mathcal{R}(A)^\perp = \mathcal{N}(A^T),$$

2. Similarly, we can prove $\mathcal{R}(A^T)^\perp = \mathcal{N}(A)$,

□

1.2 Calculus Preliminaries

Definition 1.8. (*Taylor formula for one variable*) Let $f(x)$ to be n -th differentiable at x_0 , then there exists a neighborhood $B(x_0, \epsilon)$, $\forall x \in B(x_0, \epsilon)$, s.t.

$$f(x) = f(x_0) + f'(x_0)\Delta x + \frac{f''(x_0)}{2!}\Delta x^2 + \cdots + \frac{f^{(n)}(x_0)}{n!}\Delta x^n + \mathcal{O}(\Delta x^{n+1}). \quad (17)$$

Definition 1.9. (*Taylor formula for two variables*) Let $f(x, y) \in C^{k+1}(B((x_0, y_0), \epsilon))$, then $\forall (x_0 + \Delta x, y_0 + \Delta y) \in B((x_0, y_0), \epsilon)$,

$$\begin{aligned} f(x_0 + \Delta x, y_0 + \Delta y) &= f(x_0, y_0) + \left(\Delta x \frac{\partial}{\partial x} + \Delta y \frac{\partial}{\partial y} \right) f(x_0, y_0) \\ &+ \frac{1}{2!} \left(\Delta x \frac{\partial}{\partial x} + \Delta y \frac{\partial}{\partial y} \right)^2 f(x_0, y_0) + \cdots \\ &+ \frac{1}{k!} \left(\Delta x \frac{\partial}{\partial x} + \Delta y \frac{\partial}{\partial y} \right)^k f(x_0, y_0) + \mathcal{R}_k \end{aligned} \quad (18)$$

where

$$\mathcal{R}_k = \frac{1}{(k+1)!} \left(\Delta x \frac{\partial}{\partial x} + \Delta y \frac{\partial}{\partial y} \right)^{k+1} f(x_0, y_0) f(x_0 + \theta \Delta x, y_0 + \theta \Delta y), \quad \theta \in (0, 1).$$

Definition 1.10. (*Commonly used Taylor series*)

$$\frac{1}{1-x} = \sum_{n=0}^{\infty} x^n = 1 + x + x^2 + x^3 + x^4 + \cdots, \quad x \in (-1, 1), \quad (19)$$

$$e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!} = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} + \cdots, \quad x \in \mathbb{R}, \quad (20)$$

$$\sin(x) = \sum_{n=0}^{\infty} (-1)^n \frac{x^{2n+1}}{(2n+1)!} = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \frac{x^9}{9!} - \cdots, \quad x \in \mathbb{R}, \quad (21)$$

$$\cos(x) = \sum_{n=0}^{\infty} (-1)^n \frac{x^{2n}}{(2n)!} = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \frac{x^8}{8!} - \cdots, \quad x \in \mathbb{R}, \quad (22)$$

$$\ln(1+x) = \sum_{n=0}^{\infty} (-1)^n \frac{x^{n+1}}{n+1} = x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \cdots, \quad x \in (-1, 1). \quad (23)$$

Definition 1.11. (*Cauchy's Inequality*)

$$ab \leq \frac{a^2}{2} + \frac{b^2}{2}, \quad \text{for all } a, b \in \mathbb{R}. \quad (24)$$

Definition 1.12. (Cauchy's Inequality with ϵ)

$$ab \leq \epsilon a^2 + \frac{b^2}{4\epsilon}, \quad \text{for all } a, b > 0, \epsilon > 0. \quad (25)$$

Definition 1.13. (Young's Inequality) Let $1 < p, q < \infty, \frac{1}{p} + \frac{1}{q} = 1$. Then

$$ab \leq \frac{a^p}{p} + \frac{b^q}{q} \quad \text{for all } a, b > 0. \quad (26)$$

Definition 1.14. (Young's Inequality with ϵ)

$$ab \leq \epsilon a^p + C(\epsilon)b^q, \quad \text{for all } a, b > 0, \epsilon > 0, \quad (27)$$

Where, $C(\epsilon) = (\epsilon p)^{-p/q} q^{-1}$.

Definition 1.15. (Hölder's Inequality) Let $1 < p, q < \infty, \frac{1}{p} + \frac{1}{q} = 1$. Then if $u \in L^p(U), v \in L^q(U)$, we have

$$\int_U |uv| dx \leq \left(\int_U |u|^p dx \right)^{1/p} \left(\int_U |v|^q dx \right)^{1/q} = \|u\|_{L^p(U)} \|v\|_{L^q(U)}. \quad (28)$$

Discrete Version:

$$\left| \sum_{k=1}^n a_k b_k \right| \leq \left(\sum_{k=1}^n |a_k|^p \right)^{1/p} \left(\sum_{k=1}^n |b_k|^q \right)^{1/q}. \quad (29)$$

General Version: Let $1 < p_1, \dots, p_n < \infty, \frac{1}{p_1} + \dots + \frac{1}{p_n} = 1$. Then if $u_k \in L^{p_k}(U)$, we have

$$\int_U |u_1 \cdots u_n| dx \leq \prod_{k=1}^n \|u_k\|_{L^{p_k}(U)}. \quad (30)$$

Definition 1.16. (Cauchy-Schwarz's Inequality) Let $1 \leq p < \infty$ and $u, v \in L^p(U)$. Then

$$|uv|^2 \leq \|u\|_{L^2(U)} \|v\|_{L^2(U)}. \quad (31)$$

Discrete Version:

$$\left| \sum_{i=1}^n x_i y_i \right|^2 \leq \sum_{i=1}^n |x_i|^2 \sum_{i=1}^n |y_i|^2. \quad (32)$$

Definition 1.17. (*Minkowski's Inequality*) Let $1 \leq p < \infty$ and $u, v \in L^p(U)$. Then

$$\|u + v\|_{L^p(U)} = \|u\|_{L^p(U)} + \|v\|_{L^p(U)}. \quad (33)$$

Discrete Version:

$$\left(\sum_{k=1}^n |a_k + b_k|^p \right)^{1/p} \leq \left(\sum_{k=1}^n |a_k|^p \right)^{1/p} + \left(\sum_{k=1}^n |b_k|^p \right)^{1/p}. \quad (34)$$

Definition 1.18. (*Gronwall's Inequality*)

Differential Version: Let $\eta(\cdot)$ be a *nonnegative, absolutely continuous function* on $[0, T]$, which satisfies for a.e t the differential inequality

$$\eta'(t) \leq \phi(t)\eta(t) + \psi(t), \quad (35)$$

where $\phi(t)$ and $\psi(t)$ are *nonnegative, summable functions* on $[0, T]$. Then

$$\eta(t) \leq e^{\int_0^t \phi(s)ds} \left[\eta(0) + \int_0^t \psi(s)ds \right], \forall 0 \leq t \leq T. \quad (36)$$

In particular, if

$$\eta' \leq \phi\eta, \text{ on } [0, T] \text{ and } \eta(0) = 0, \quad (37)$$

$$\eta(t) = 0, \forall 0 \leq t \leq T. \quad (38)$$

Integral Version: Let $\xi(\cdot)$ be a *nonnegative, summable function* on $[0, T]$, which satisfies for a.e t the integral inequality

$$\xi(t) \leq C_1 \int_0^t \xi(s)ds + C_2, \quad (39)$$

where $C_1, C_2 \geq 0$. Then

$$\xi(t) \leq C_2 \left(1 + C_1 t e^{C_1 t} \right), \forall a.e. \ 0 \leq t \leq T. \quad (40)$$

In particular, if

$$\xi(t) \leq C_1 \int_0^t \xi(s)ds, \forall a.e. \ 0 \leq t \leq T, \quad (41)$$

$$\xi(t) = 0, a.e. \quad (42)$$

Discrete Version: If

$$(1 + \gamma)a_{n+1} \leq a_n + \beta f_n, \ \alpha, \gamma \in \mathbb{R} \ \gamma \neq -1, \quad (43)$$

then,

$$a_{n+1} \leq \frac{a_0}{(1 + \gamma)^{n+1}} + \beta \sum_{k=0}^n \frac{f_k}{(1 + \gamma)^{n+k-1}}. \quad (44)$$

1.3 Norms' Preliminaries

1.3.1 Vector Norms

Definition 1.19. (Vector Norms) A vector norm is a function $\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfying the following conditions for all $x, y \in \mathbb{R}^n$ and $\alpha \in \mathbb{R}$

1. *nonnegative* : $\|x\| \geq 0$, ($\|x\| = 0 \Leftrightarrow x = 0$),
2. *homogeneity* : $\|\alpha x\| = |\alpha| \|x\|$,
3. *triangle inequality* : $\|x + y\| \leq \|x\| + \|y\|$, $\forall x, y \in \mathbb{R}^n$,

Definition 1.20. For $x \in \mathbb{R}^n$, some of the most frequently used vector norms are

1. *1-norm* : $\|x\|_1 = \sum_{i=1}^n |x_i|$,
2. *2-norm* : $\|x\|_2 = \sqrt{\sum_{i=1}^n |x_i|^2}$,
3. *∞ -norm* : $\|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$,
4. *p-norm* : $\|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}$.

Corollary 1.3. For all $x \in \mathbb{R}^n$,

$$\|x\|_2 \leq \|x\|_1 \leq \sqrt{n} \|x\|_2, \quad (45)$$

$$\|x\|_\infty \leq \|x\|_2 \leq \sqrt{n} \|x\|_\infty, \quad (46)$$

$$\frac{1}{\sqrt{n}} \|x\|_1 \leq \|x\|_2 \leq \sqrt{n} \|x\|_1, \quad (47)$$

$$\|x\|_\infty \leq \|x\|_1 \leq \sqrt{n} \|x\|_\infty. \quad (48)$$

Theorem 1.22. (vector 2-norm invariance) Vector 2-norm is invariant under the orthogonal transformation, i.e., if Q is an $n \times n$ orthogonal matrix, then

$$\|Qx\|_2 = \|x\|_2, \quad \forall x \in \mathbb{R}^n \quad (49)$$

Proof.

$$\|Qx\|_2^2 = (Qx)^T (Qx) = x^T Q^T Q x = x^T x = \|x\|_2^2.$$

□

1.3.2 Matrix Norms

Definition 1.21. (Matrix Norms) A matrix norm is a function $\|\cdot\| : \mathbb{R}^{m \times n} \mapsto \mathbb{R}$ satisfying the following conditions for all $A, B \in \mathbb{R}^{m \times n}$ and $\alpha \in \mathbb{R}$

1. *nonnegative* : $\|x\| \geq 0$, ($\|x\| = 0 \Leftrightarrow x = 0$),
2. *homogeneity* : $\|\alpha x\| = |\alpha| \|x\|$,
3. *triangle inequality* : $\|x + y\| \leq \|x\| + \|y\|$, $\forall x, y \in \mathbb{R}^n$,

Definition 1.22. For $A \in \mathbb{R}^{m \times n}$, some of the most frequently matrix vector norms are

1. *F-norm* : $\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}$,
3. *∞ -norm* : $\|A\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}|$,
2. *1-norm* : $\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}|$,
4. *induced-norm* : $\|A\|_p = \sup_{x \in \mathbb{R}^n, x \neq 0} \frac{\|Ax\|_p}{\|x\|_p}$.

Corollary 1.4. For all $A \in \mathbb{C}^{n \times n}$,

$$\|A\|_2 \leq \|A\|_F \leq \sqrt{n} \|A\|_2, \quad (50)$$

$$\frac{1}{\sqrt{n}} \|A\|_2 \leq \|A\|_\infty \leq \sqrt{n} \|A\|_2, \quad (51)$$

$$\frac{1}{\sqrt{n}} \|A\|_\infty \leq \|A\|_2 \leq \sqrt{n} \|A\|_\infty, \quad (52)$$

$$\frac{1}{\sqrt{n}} \|A\|_1 \leq \|A\|_2 \leq \sqrt{n} \|A\|_1. \quad (53)$$

Corollary 1.5. For all $A \in \mathbb{C}^{n \times n}$, then $\|A\|_2 \leq \sqrt{\|A\|_1 \|A\|_\infty}$.

Proof.

$$\|A\|_2^2 = \kappa(A) \leq \|A\|_1 \|A^{-1}\|_1 = \|A\|_1 \|A\|_\infty.$$

$\kappa(A)$: condition number □

Theorem 1.23. (Matrix 2-norm and Frobenius invariance) (Matrix 2-norm and Frobenius are invariant under the orthogonal transformation, i.e., if Q is an $n \times n$ orthogonal matrix, then

$$\|QA\|_2 = \|A\|_2, \quad \forall A \in \mathbb{R}^{n \times n}, \quad (54)$$

$$\|QA\|_F = \|A\|_F, \quad \forall A \in \mathbb{R}^{n \times n} \quad (55)$$

Theorem 1.24. Suppose that $A \in \mathbb{R}^{n \times n}$. If $\|A\| < 1$, then $(I - A)$ is nonsingular and

$$(I - A)^{-1} = \sum_{k=0}^{\infty} A^k \quad (56)$$

with

$$\|(I - A)^{-1}\| \leq \frac{1}{1 - \|A\|}. \quad (57)$$

Lemma 1.1. Suppose that $A \in \mathbb{R}^{n \times n}$. If $(I - A)$ is singular, then $\|A\| \geq 1$.

9

10

2 Direct Method

2.1 For square or rectangular matrices $A \in \mathbb{C}^{m,n}, m \geq n$

2.1.1 Singular Value Decomposition

2.1.2 QR Decomposition

2.2 For square matrices $A \in \mathbb{C}^{n,n}$

2.2.1 LU Decomposition

2.2.2 Cholesky Decomposition

3 Iterative Method

3.1 General Iterative Scheme

3.2 Jacobi Method

3.3 Gauss-Seidel Method

3.4 Richardson Method

3.5 Successive Over Relaxation (SOR) Method

3.6 Minimal Correction Method

3.7 Steepest Descent Method

3.8 Conjugate Gradients Method

4 Eigenvalue Problems

5 Euler Method

In this section, we focus on

$$\begin{cases} y' = f(t, y), \\ y(t_0) = y_0. \end{cases}$$

Where f is [Lipschitz continuous](#) w.r.t. the second variable, i.e

$$|f(t, x) - f(t, y)| \leq \lambda |x - y|, \quad \lambda > 0. \quad (58)$$

In the following, We will let $y(t_n)$ to be the numerical approximation of y_n and $e_n = y_n - y(t_n)$ to be the error.

Definition 5.1. (*Order of the Method*) A time stepping scheme

$$y_{n+1} = \Phi(h, y_0, y_1, \dots, y_n) \quad (59)$$

is of order of $p \geq 1$, if

$$y_{n+1} - \Phi(h, y_0, y_1, \dots, y_n) = \mathcal{O}(h^{p+1}). \quad (60)$$

Definition 5.2. (*Convergence of the Method*) A time stepping scheme

$$y_{n+1} = \Phi(h, y_0, y_1, \dots, y_n) \quad (61)$$

is convergent, if

$$\lim_{h \rightarrow 0} \max_n \|y(t_n) - y_n\| = 0. \quad (62)$$

5.1 Euler's method

Definition 5.3. (*Forward Euler Method^a*)

$$y_{n+1} = y_n + hf(t_n, y_n), \quad n = 0, 1, 2, \dots \quad (63)$$

^aForward Euler Method is explicit.

Theorem 5.1. (*Forward Euler Method is of order 1^a*) Forward Euler Method

$$y(t_{n+1}) = y(t_n) + hf(t_n, y(t_n)), \quad (64)$$

is of order 1.

^aYou can also use multi-step theorem to derive it.

Proof. By the Taylor expansion,

$$y(t_{n+1}) = y(t_n) + hy'(t_n) + \mathcal{O}(h^2). \quad (65)$$

So,

$$\begin{aligned} y(t_{n+1}) - y(t_n) - hf(t_n, y(t_n)) &= y(t_n) + hy'(t_n) + \mathcal{O}(h^2) - y(t_n) - hf(t_n, y(t_n)) \\ &= y(t_n) + hy'(t_n) + \mathcal{O}(h^2) - y(t_n) - hy'(t_n) \\ &= \mathcal{O}(h^2). \end{aligned} \quad (66)$$

Therefore, Forward Euler Method (5.3) is order of 1. □

Theorem 5.2. (*The convergence of Forward Euler Method*) Forward Euler Method

$$y(t_{n+1}) = y(t_n) + hf(t_n, y(t_n)), \quad (67)$$

is convergent.

Proof. From (66), we get

$$y(t_{n+1}) = y(t_n) + hf(t_n, y(t_n)) + \mathcal{O}(h^2), \quad (68)$$

Subtracting (68) from (63), we get

$$e_{n+1} = e_n + h[f(t_n, y_n) - f(t_n, y(t_n))] + ch^2. \quad (69)$$

Since f is lipschitz continuous w.r.t. the second variable, then

$$|f(t_n, y_n) - f(t_n, y(t_n))| \leq \lambda |y_n - y(t_n)|, \quad \lambda > 0. \quad (70)$$

Therefore,

$$\begin{aligned} \|e_{n+1}\| &\leq \|e_n\| + h\lambda \|e_n\| + ch^2 \\ &= (1 + h\lambda) \|e_n\| + ch^2. \end{aligned} \quad (71)$$

Claim:[?]

$$\|e_n\| \leq \frac{c}{\lambda} h[(1 + h\lambda)^n - 1], n = 0, 1, \dots \quad (72)$$

Proof for Claim (72): The proof is by induction on n .

1. when $n = 0$, $e_n = 0$, hence $\|e_n\| \leq \frac{c}{\lambda} h[(1 + h\lambda)^n - 1]$,

2. Induction assumption:

$$\|e_n\| \leq \frac{c}{\lambda} h[(1 + h\lambda)^n - 1]$$

3. Induction steps:

$$\|e_{n+1}\| \leq (1 + h\lambda) \|e_n\| + ch^2 \quad (73)$$

$$\leq (1 + h\lambda) \frac{c}{\lambda} h[(1 + h\lambda)^n - 1] + ch^2 \quad (74)$$

$$= \frac{c}{\lambda} h[(1 + h\lambda)^{n+1} - 1]. \quad (75)$$

So, from the claim (72), we get $\|e_n\| \rightarrow 0$, when $h \rightarrow 0$. Therefore **Forward Euler Method is convergent**. \square

Definition 5.4. (*Backward Euler Methods^a*)

$$y_{n+1} = y_n + hf(t_{n+1}, y_{n+1}), \quad n = 0, 1, 2, \dots. \quad (76)$$

^aBackward Euler Method is implicit.

Theorem 5.3. (*backward Euler Method is of order 1*^a) Backward Euler Method

$$y(t_{n+1}) = y(t_n) + hf(t_{n+1}, y(t_{n+1})), \quad (77)$$

is of order 1 .

^aYou can also use multi-step theorem to derive it.

Proof. By the Taylor expansion,

$$y(t_{n+1}) = y(t_n) + hy'(t_n) + \mathcal{O}(h^2) \quad (78)$$

$$y'(t_{n+1}) = y'(t_n) + \mathcal{O}(h). \quad (79)$$

So,

$$\begin{aligned} & y(t_{n+1}) - y(t_n) - hf(t_{n+1}, y(t_{n+1})) \\ &= y(t_{n+1}) - y(t_n) + hy'(t_{n+1}) \\ &= y(t_n) + hy'(t_n) + \mathcal{O}(h^2) - y(t_n) - h[y'(t_n) + \mathcal{O}(h)] \\ &= \mathcal{O}(h^2). \end{aligned} \quad (80)$$

Therefore, Backward Euler Method (5.4) is order of 1 . □

Theorem 5.4. (*The convergence of Backward Euler Method*) Backward Euler Method

$$y(t_{n+1}) = y(t_n) + hf(t_{n+1}, y(t_{n+1})), \quad (81)$$

is convergent.

Proof. From (80), we get

$$y(t_{n+1}) = y(t_n) + hf(t_{n+1}, y(t_{n+1})) + \mathcal{O}(h^2), \quad (82)$$

Subtracting (82) from (76), we get

$$e_{n+1} = e_n + h[f(t_{n+1}, y_{n+1}) - f(t_{n+1}, y(t_{n+1}))] + \mathcal{O}(h^2). \quad (83)$$

Since f is lipschitz continuous w.r.t. the second variable, then

$$|f(t_{n+1}, y_{n+1}) - f(t_{n+1}, y(t_{n+1}))| \leq \lambda |y_{n+1} - y(t_{n+1})|, \quad \lambda > 0. \quad (84)$$

Therefore,

$$\|e_{n+1}\| \leq \|e_n\| + h\lambda \|e_{n+1}\| + \mathcal{O}(h^2). \quad (85)$$

So,

$$(1 - h\lambda)\|e_{n+1}\| \leq \|e_n\| + \mathcal{O}(h^2). \quad (86)$$

So, by the [Discrete Gronwall's Inequality](#), we have

$$\begin{aligned}
 \|e_{n+1}\| &\leq \frac{\|e_0\|}{(1-h\lambda)^{n+1}} + c \sum_{k=0}^n \frac{h^2}{(1-h\lambda)^{n+k-1}} \\
 &= c \sum_{k=0}^n \frac{h^2}{(1-h\lambda)^{n+k-1}} \\
 &\leq ch^2(1+h\lambda)^{(nh)/h\lambda}(1-h\lambda \rightarrow 1+h\lambda) \\
 &\leq che^T T.
 \end{aligned} \tag{87}$$

So, from the claim (87), we get $\|e_n\| \rightarrow 0$, when $h \rightarrow 0$. Therefore [Forward Euler Method is convergent](#). \square

5.2 Trapezoidal Method

Definition 5.5. (*Trapezoidal Method^a*)

$$y_{n+1} = y_n + \frac{1}{2}h[f(t_n, y_n) + f(t_{n+1}, y_{n+1})], \quad n = 0, 1, 2, \dots \tag{88}$$

^aTrapezoidal Method Method is a combination of Forward Euler Method and Backward Euler Method.

Theorem 5.5. (*Trapezoidal Method is of order 2^a*) Trapezoidal Method

$$y(t_{n+1}) = y(t_n) + \frac{1}{2}h[f(t_n, y(t_n)) + f(t_{n+1}, y(t_{n+1}))], \tag{89}$$

is of order 2.

^aYou can also use multi-step theorem to derive it.

Proof. By the Taylor expansion,

$$y(t_{n+1}) = y(t_n) + hy'(t_n) + \frac{1}{2!}h^2y''(t_n) + \mathcal{O}(h^3) \tag{90}$$

$$y'(t_{n+1}) = y'(t_n) + hy''(t_n) + \mathcal{O}(h^2). \tag{91}$$

So,

$$\begin{aligned}
 &y(t_{n+1}) - y(t_n) + \frac{1}{2}h[f(t_n, y(t_n)) + f(t_{n+1}, y(t_{n+1}))] \\
 &= y(t_{n+1}) - y(t_n) + \frac{1}{2}h[y'(t_n) + y'(t_{n+1})] \\
 &= y(t_n) + hy'(t_n) + \frac{1}{2!}h^2y''(t_n) + \mathcal{O}(h^3) - y(t_n) + \frac{1}{2}h[y'(t_n) + y'(t_n) + hy''(t_n) + \mathcal{O}(h^2)] \\
 &= \mathcal{O}(h^3).
 \end{aligned} \tag{92}$$

Therefore, Trapezoidal Method (5.5) is order of 2. \square

Theorem 5.6. (*The convergence of Trapezoidal Method*) Trapezoidal Method

$$y(t_{n+1}) = y(t_n) + \frac{1}{2}h[f(t_n, y(t_n)) + f(t_{n+1}, y(t_{n+1}))], \quad (93)$$

is convergent.

Proof. From (92), we get

$$y(t_{n+1}) = y(t_n) + \frac{1}{2}h[f(t_n, y(t_n)) + f(t_{n+1}, y(t_{n+1}))] + \mathcal{O}(h^3), \quad (94)$$

Subtracting (94) from (88), we get

$$e_{n+1} = e_n + \frac{1}{2}h[f(t_n, y_n) - f(t_n, y(t_n)) + f(t_{n+1}, y_{n+1}) - f(t_{n+1}, y(t_{n+1}))] + ch^3. \quad (95)$$

Since f is lipschitz continuous w.r.t. the second variable, then

$$|f(t_n, y_n) - f(t_n, y(t_n))| \leq \lambda|y_n - y(t_n)|, \quad \lambda > 0, \quad (96)$$

$$|f(t_{n+1}, y_{n+1}) - f(t_{n+1}, y(t_{n+1}))| \leq \lambda|y_{n+1} - y(t_{n+1})|, \quad \lambda > 0. \quad (97)$$

Therefore,

$$\|e_{n+1}\| \leq \|e_n\| + \frac{1}{2}h\lambda(\|e_n\| + \|e_{n+1}\|) + ch^3. \quad (98)$$

So,

$$(1 - \frac{1}{2}h\lambda)\|e_{n+1}\| \leq (1 + \frac{1}{2}h\lambda)\|e_n\| + ch^3. \quad (99)$$

Claim:[?]

$$\|e_n\| \leq \frac{c}{\lambda}h^2 \left[\left(\frac{1 + \frac{1}{2}h\lambda}{1 - \frac{1}{2}h\lambda} \right)^n - 1 \right], n = 0, 1, \dots \quad (100)$$

Proof for Claim (100): The proof is by induction on n .

Then, we can make h small enough to such that $0 < h\lambda < 2$, then

$$\frac{1 + \frac{1}{2}h\lambda}{1 - \frac{1}{2}h\lambda} = 1 + \frac{1}{1 - \frac{1}{2}h\lambda} \leq \sum_{\ell=0}^{\infty} \frac{1}{\ell!} \left(\frac{h\lambda}{1 - \frac{1}{2}h\lambda} \right)^{\ell} = \exp\left(\frac{h\lambda}{1 - \frac{1}{2}h\lambda} \right).$$

Therefore,

$$\|e_n\| \leq \frac{c}{\lambda}h^2 \left[\left(\frac{1 + \frac{1}{2}h\lambda}{1 - \frac{1}{2}h\lambda} \right)^n - 1 \right] \leq \frac{c}{\lambda}h^2 \left(\frac{1 + \frac{1}{2}h\lambda}{1 - \frac{1}{2}h\lambda} \right)^n \leq \frac{c}{\lambda}h^2 \exp\left(\frac{nh\lambda}{1 - \frac{1}{2}h\lambda} \right). \quad (101)$$

This bound of true for every negative integer n such that $nh < T$. Therefore,

$$\|e_n\| \leq \frac{c}{\lambda}h^2 \exp\left(\frac{nh\lambda}{1 - \frac{1}{2}h\lambda} \right) \leq \frac{c}{\lambda}h^2 \exp\left(\frac{T\lambda}{1 - \frac{1}{2}h\lambda} \right). \quad (102)$$

So, from the claim (102), we get $\|e_n\| \rightarrow 0$, when $h \rightarrow 0$. Therefore **Forward Euler Method is convergent**. \square

5.3 Theta Method

Definition 5.6. (*Theta Method^a*)

$$y_{n+1} = y_n + \frac{1}{2}h[\theta f(t_n, y_n) + (1 - \theta)f(t_{n+1}, y_{n+1})], \quad n = 0, 1, 2, \dots \quad (103)$$

^aTheta Method is a general form of Forward Euler Method ($\theta = 1$), Backward Euler Method ($\theta = 0$) and Trapezoidal Method ($\theta = \frac{1}{2}$).

5.4 Midpoint Rule Method

Definition 5.7. (*Midpoint Rule Method*)

$$y_{n+1} = y_n + hf\left(t_n + \frac{1}{2}h, \frac{1}{2}(y_n + y_{n+1})\right). \quad (104)$$

Theorem 5.7. (*Midpoint Rule Method is of order 2*) Midpoint Rule Method

$$y(t_{n+1}) = y(t_n) + hf\left(t_n + \frac{1}{2}h, \frac{1}{2}(y(t_n) + y(t_{n+1}))\right). \quad (105)$$

is of order 2.

Proof. By the Taylor expansion,

$$y(t_{n+1}) = y(t_n) + hy'(t_n) + \frac{1}{2!}h^2y''(t_n) + \mathcal{O}(h^3) \quad (106)$$

$$f(x_0 + \Delta x, y_0 + \Delta y) = f(x_0, y_0) + \left(\Delta x \frac{\partial}{\partial x} + \Delta y \frac{\partial}{\partial y}\right)f(x_0, y_0) + \mathcal{O}(h^2). \quad (107)$$

And chain rule

$$y'' = f'(t, \mathbf{y}) = \frac{\partial f(t, \mathbf{y})}{\partial t} + \frac{\partial f(t, \mathbf{y})}{\partial \mathbf{y}} f(t, \mathbf{y}). \quad (108)$$

So,

$$\begin{aligned} & y(t_{n+1}) - y(t_n) + hf\left(t_n + \frac{1}{2}h, \frac{1}{2}(y(t_n) + y(t_{n+1}))\right) \\ &= y(t_n) + hy'(t_n) + \frac{1}{2!}h^2y''(t_n) + \mathcal{O}(h^3) - y(t_n) \\ &= h\left(f(t_n, y_n) + (t_n + \frac{1}{2}h - t_n)\frac{\partial f(t_n, y_n)}{\partial t} + (\frac{1}{2}(y(t_n) + y(t_{n+1})) - y_n)\frac{\partial f(t_n, y_n)}{\partial y} + \mathcal{O}(h^2)\right) \\ &= hy'(t_n) + \frac{1}{2!}h^2y''(t_n) + \mathcal{O}(h^3) \\ &= \left(hf(t_n, y_n) + \frac{1}{2}h^2\frac{\partial f(t_n, y_n)}{\partial t} + \frac{1}{2}h^2\frac{\partial f(t_n, y_n)}{\partial y} + \mathcal{O}(h^3)\right) \end{aligned}$$

$$\begin{aligned}
&= hy'(t_n) + \frac{1}{2!}h^2 \left(\frac{\partial f(t_n, y_n)}{\partial t} + \frac{\partial f(t_n, y_n)}{\partial y} y'(t_n) \right) \\
&- \left(hy'(t_n) + \frac{1}{2}h^2 \frac{\partial f(t_n, y_n)}{\partial t} + \frac{1}{2}h^2 \frac{\partial f(t_n, y_n)}{\partial y} + \mathcal{O}(h^3) \right) \\
&= \mathcal{O}(h^3).
\end{aligned}$$

Therefore, Midpoint Rule Method (5.5) is order of 2. □

Theorem 5.8. (*The convergence of Midpoint Rule Method*) Midpoint Rule Method

$$y(t_{n+1}) = y(t_n) + hf \left(t_n + \frac{1}{2}h, \frac{1}{2}(y(t_n) + y(t_{n+1})) \right), \quad (109)$$

is convergent.

Proof. From (109), we get

$$y(t_{n+1}) = y(t_n) + hf \left(t_n + \frac{1}{2}h, \frac{1}{2}(y(t_n) + y(t_{n+1})) \right) + \mathcal{O}(h^3), \quad (110)$$

Subtracting (110) from (104), we get

$$e_{n+1} = e_n + h \left[f \left(t_n + \frac{1}{2}h, \frac{1}{2}(y(t_n) + y(t_{n+1})) \right) - f \left(t_n + \frac{1}{2}h, \frac{1}{2}(y(t_n) + y(t_{n+1})) \right) \right] + ch^3. \quad (111)$$

Since f is lipschitz continuous w.r.t. the second variable, then

$$\begin{aligned}
&\left| f \left(t_n + \frac{1}{2}h, \frac{1}{2}(y(t_n) + y(t_{n+1})) \right) - f \left(t_n + \frac{1}{2}h, \frac{1}{2}(y(t_n) + y(t_{n+1})) \right) \right| \\
&\leq \frac{1}{2}\lambda |y_n - y(t_n) + y_{n+1} - y(t_{n+1})|, \quad \lambda > 0.
\end{aligned} \quad (112)$$

Therefore,

$$\|e_{n+1}\| \leq \|e_n\| + \frac{1}{2}h\lambda(\|e_n\| + \|e_{n+1}\|) + ch^3. \quad (113)$$

So,

$$(1 - \frac{1}{2}h\lambda)\|e_{n+1}\| \leq (1 + \frac{1}{2}h\lambda)\|e_n\| + ch^3. \quad (114)$$

Claim:[?]

$$\|e_n\| \leq \frac{c}{\lambda} h^2 \left[\left(\frac{1 + \frac{1}{2}h\lambda}{1 - \frac{1}{2}h\lambda} \right)^n - 1 \right], n = 0, 1, \dots \quad (115)$$

Proof for Claim (115): The proof is by induction on n .

Then, we can make h small enough to such that $0 < h\lambda < 2$, then

$$\frac{1 + \frac{1}{2}h\lambda}{1 - \frac{1}{2}h\lambda} = 1 + \frac{1}{1 - \frac{1}{2}h\lambda} \leq \sum_{\ell=0}^{\infty} \frac{1}{\ell!} \left(\frac{h\lambda}{1 - \frac{1}{2}h\lambda} \right)^{\ell} = \exp\left(\frac{h\lambda}{1 - \frac{1}{2}h\lambda} \right).$$

Therefore,

$$\|e_n\| \leq \frac{c}{\lambda} h^2 \left[\left(\frac{1 + \frac{1}{2}h\lambda}{1 - \frac{1}{2}h\lambda} \right)^n - 1 \right] \leq \frac{c}{\lambda} h^2 \left(\frac{1 + \frac{1}{2}h\lambda}{1 - \frac{1}{2}h\lambda} \right)^n \leq \frac{c}{\lambda} h^2 \exp\left(\frac{nh\lambda}{1 - \frac{1}{2}h\lambda} \right). \quad (116)$$

This bound of true for every negative integer n such that $nh < T$. Therefore,

$$\|e_n\| \leq \frac{c}{\lambda} h^2 \exp\left(\frac{nh\lambda}{1 - \frac{1}{2}h\lambda} \right) \leq \frac{c}{\lambda} h^2 \exp\left(\frac{T\lambda}{1 - \frac{1}{2}h\lambda} \right). \quad (117)$$

So, from the claim (117), we get $\|e_n\| \rightarrow 0$, when $h \rightarrow 0$. Therefore **Midpoint Rule Method is convergent**. \square

6 Multistep Method

6.1 The Adams Method

Definition 6.1. (*s-step Adams-bashforth*)

$$y_{n+s} = y_{n+s-1} + h \sum_{m=0}^{s-1} b_m f(t_{n+m}, y_{n+m}), \quad (118)$$

where

$$b_m = h^{-1} \int_{t_{n+s-1}}^{t_{n+s}} p_m(\tau) d\tau = h^{-1} \int_0^h p_m(t_{n+s-1} + \tau) d\tau \quad n = 0, 1, 2, \dots.$$

$$p_m(t) = \prod_{l=0, l \neq m}^{s-1} \frac{t - t_{n+l}}{t_{n+m} - t_{n+l}}, \quad \text{Lagrange interpolation polynomials}.$$

(1-step Adams-bashforth)

$$y_{n+1} = y_n + h f(t_n, y_n),$$

(2-step Adams-bashforth)

$$y_{n+2} = y_{n+1} + h \left[\frac{3}{2} f(t_{n+1}, y_{n+1}) - \frac{1}{2} f(t_n, y_n) \right],$$

(3-step Adams-bashforth)

$$y_{n+3} = y_{n+2} + h \left[\frac{23}{12} f(t_{n+2}, y_{n+2}) - \frac{4}{3} f(t_{n+1}, y_{n+1}) + \frac{5}{12} f(t_n, y_n) \right].$$

6.2 The Order and Convergence of Multistep Methods

Definition 6.2. (*General s-step Method*) The general s-step Method ^a can be written as

$$\sum_{m=0}^s a_m y_{n+m} = h \sum_{m=0}^s b_m f(t_{n+m}, y_{n+m}). \quad (119)$$

Where $a_m, b_m, m = 0, \dots, s$, are given constants, independent of h, n and original equation.

^aif $b_s = 0$ the method is explicit; otherwise it is implicit.

Theorem 6.1. (*s-step method convergent order*) The multistep method (119) is of order $p \geq 1$ if and only if there exists $c \neq 0$ s.t.

$$\rho(w) - \sigma(w) \ln w^a = c(w-1)^{p+1} + \mathcal{O}(|w-1|^{p+2}), \quad w \rightarrow 1. \quad (120)$$

Where,

$$\rho(w) := \sum_{m=0}^s a_m w^m \quad \text{and} \quad \sigma(w) := \sum_{m=0}^s b_m w^m. \quad (121)$$

^aLet $w = \xi + 1$, then $\ln(1 + \xi) = \sum_{n=0}^{\infty} (-1)^n \frac{\xi^{n+1}}{n+1} = \xi - \frac{\xi^2}{2} + \frac{\xi^3}{3} - \frac{\xi^4}{4} + \dots + (-1)^n \frac{\xi^{n+1}}{n+1} + \dots, \xi \in (-1, 1)$.

Theorem 6.2. (*s-step method convergent order*) The multistep method (119) is of order $p \geq 1$ if and only if

1. $\sum_{m=0}^s a_m = 0$, (i.e. $\rho(1) = 0$),
2. $\sum_{m=0}^s m^k a_m = k \sum_{m=0}^s m^{k-1} b_m, k = 1, 2, \dots, p$,
3. $\sum_{m=0}^s m^{p+1} a_m \neq (p+1) \sum_{m=0}^s m^p b_m$.

Where,

$$\rho(w) := \sum_{m=0}^s a_m w^m \quad \text{and} \quad \sigma(w) := \sum_{m=0}^s b_m w^m. \quad (122)$$

Lemma 6.1. (*Root Condition*) If the roots $|\lambda_i| \leq 1$ for each $i = 1, \dots, m$ and all roots with value 1 are simple root then the difference method is said to satisfy the root condition.

Theorem 6.3. (*The Dahlquist equivalence theorem*) The multistep method (119) is convergent if and only if

1. *consistency*: multistep method (119) is order of $p \geq 1$,
2. *stability*: the polynomial $\rho(w)$ satisfies the root condition.

6.3 Method of A-stable verification for Multistep Methods

Theorem 6.4. *Explicit Multistep Methods can not be A-stable.*

Theorem 6.5. *(Dahlquist second barrier) The highest order of an A-stable multistep method is 2.*

7 Runge-Kutta Methods

7.1 Quadrature Formulas

Definition 7.1. *(The Quadrature) The Quadrature is the procedure of replacing an integral with a finite sum.*

Definition 7.2. *(The Quadrature Formula) Let w be a nonnegative function in (a,b) s.t.*

$$0 < \int_a^b w(\tau) d\tau < \infty, \quad \left| \int_a^b \tau^j w(\tau) d\tau \right| < \infty, j = 1, 2, \dots$$

Then, the quadrature formula is as following

$$\int_a^b f(\tau) w(\tau) d\tau \approx \sum_j^n b_j f(c_j). \quad (123)$$

Remark 7.1. *The quadrature formula (123) is order of p if it is exact for every $f \in \mathbb{P}_{p-1}$.*

7.2 Explicit Runge-Kutta Formulas

7.3 Implicit Runge-Kutta Formulas

7.4 Collocation Runge-Kutta Formulas

7.5 Method of A-stable verification for Runge-Kutta Method

Theorem 7.1. *Explicit Runge-Kutta Methods can not be A-stable.*

7.6 Problems

Problem 7.1. Find the order of the following quadrature formula.

$$\int_0^1 f(\tau) d\tau = \frac{1}{6}f(0) + \frac{2}{3}f\left(\frac{1}{2}\right) + \frac{1}{6}f(1), \quad \text{Simpson Rule.}$$

Solution. Since the quadrature formula (123) is order of p if it is exact for every $f \in \mathbb{P}_{p-1}$. we can chose the simplest basis $(1, \tau, \tau^2, \tau^3, \dots, \tau^{p-1})$, and the order conditions read that

$$\sum_{j=1}^p b_j c_j^m = \int_a^b \tau^m w(\tau) d\tau, \quad m = 0, 1, \dots, p-1. \quad (124)$$

Checking the order condition by the following procedure,

$$\begin{aligned} 1 &= \int_0^1 1 d\tau = \frac{1}{6}1 + \frac{2}{3}1 + \frac{1}{6}1 = 1. \\ \frac{1}{2} &= \int_0^1 \tau d\tau = \frac{1}{6}0 + \frac{2}{3}\left(\frac{1}{2}\right) + \frac{1}{6}1 = \frac{1}{2}. \\ \frac{1}{3} &= \int_0^1 \tau^2 d\tau = \frac{1}{6}0^2 + \frac{2}{3}\left(\frac{1}{2}\right)^2 + \frac{1}{6}1^2 = \frac{1}{3}. \\ \frac{1}{4} &= \int_0^1 \tau^3 d\tau = \frac{1}{6}0^3 + \frac{2}{3}\left(\frac{1}{2}\right)^3 + \frac{1}{6}1^3 = \frac{1}{4}. \\ \frac{1}{5} &= \int_0^1 \tau^4 d\tau \neq \frac{1}{6}0^4 + \frac{2}{3}\left(\frac{1}{2}\right)^4 + \frac{1}{6}1^4 = \frac{5}{24}. \end{aligned}$$

we can get the order of the Simpson rule quadrature formula is 4. ◀

Problem 7.2. Recall Simpson's quadrature rule:

$$\int_a^b f(\tau) d\tau = \frac{b-a}{6} \left[f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right] + \mathcal{O}(|b-a|^4), \quad \text{Simpson Rule.}$$

Starting from the identity

$$y(t_{n+1}) - y(t_{n-1}) = \int_{t_{n-1}}^{t_{n+1}} f(s; y(s)) ds. \quad (125)$$

use Simpson's rule to derive a 3-step method. Determine its order and whether it is convergent.

Solution. 1. The derivation of the a 3-step method

since,

$$y(t_{n+1}) - y(t_{n-1}) = \int_{t_{n-1}}^{t_{n+1}} f(s; y(s)) ds. \quad (126)$$

Then, by Simpson's quadrature rule, we have

$$y(t_{n+1}) - y(t_{n-1}) \quad (127)$$

$$= \int_{t_{n-1}}^{t_{n+1}} f(s; y(s)) ds. \quad (128)$$

$$= \frac{t_{n+1} - t_{n-1}}{6} \left[f(t_{n-1}; y(t_{n-1})) + 4f\left(\frac{t_{n-1} + t_{n+1}}{2}; y\left(\frac{t_{n-1} + t_{n+1}}{2}\right)\right) + f(t_{n+1}; y(t_{n+1})) \right] \quad (129)$$

$$= \frac{h}{3} [f(t_{n-1}; y(t_{n-1})) + 4f(t_n; y(t_n)) + f(t_{n+1}; y(t_{n+1}))]. \quad (130)$$

Therefore, the 3-step method deriving from Simpson's rule is

$$y(t_{n+1}) = y(t_{n-1}) + \frac{h}{3} [f(t_{n-1}; y(t_{n-1})) + 4f(t_n; y(t_n)) + f(t_{n+1}; y(t_{n+1}))]. \quad (131)$$

Or

$$y(t_{n+2}) - y(t_n) = \frac{h}{3} [f(t_n; y(t_n)) + 4f(t_{n+1}; y(t_{n+1})) + f(t_{n+2}; y(t_{n+2}))]. \quad (132)$$

2. **The order** For our this problem

$$\rho(w) := \sum_{m=0}^s a_m w^m = -1 + w^2 \quad \text{and} \quad \sigma(w) := \sum_{m=0}^s b_m w^m = \frac{1}{3} + \frac{4}{3}w + \frac{1}{3}w^2. \quad (133)$$

By making the substitution with $\xi = w - 1$ i.e. $w = \xi + 1$, then

$$\rho(w) := \sum_{m=0}^s a_m w^m = \xi^2 + 2\xi \quad \text{and} \quad \sigma(w) := \sum_{m=0}^s b_m w^m = \frac{1}{3}\xi^2 + 2\xi + 2. \quad (134)$$

So,

$$\begin{aligned} \rho(w) - \sigma(w) \ln(w) &= \xi^2 + 2\xi - (2 + 2\xi + \frac{1}{3}\xi^2) \left(\xi - \frac{\xi^2}{2} + \frac{\xi^3}{3} \cdots \right) \\ &= \begin{array}{ccccccc} +2\xi & +\xi^2 & & & & & \\ -2\xi & +\xi^2 & -\frac{2}{3}\xi^3 & & & & \\ & -2\xi^2 & +\xi^3 & -\frac{2}{3}\xi^4 & & & \\ & & -\frac{1}{3}\xi^3 & +\frac{1}{6}\xi^4 & -\frac{1}{9}\xi^5 & & \end{array} \\ &= -\frac{1}{2}\xi^4 + \mathcal{O}(\xi^5). \end{aligned}$$

Therefore, by the theorem

$$\rho(w) - \sigma(w) \ln(w) = -\frac{1}{2}\xi^4 + \mathcal{O}(\xi^5).$$

Hence, this scheme is order of 3.

3. The stability Since,

$$\rho(w) := \sum_{m=0}^s a_m w^m = -1 + w^2 = (w-1)(w+1). \quad (135)$$

And $w = \pm 1$ are simple root which satisfy the root condition. Therefore, this scheme is stable.

Hence, it is of order 3 and convergent. convergent

Problem 7.3. Restricting your attention to scalar autonomous $y' = f(y)$, prove that the ERK method with tableau

0				
$\frac{1}{2}$	$\frac{1}{2}$			
$\frac{1}{2}$	0	$\frac{1}{2}$		
1	0	0	1	
	$\frac{1}{6}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{6}$

is of order 4.

Solution.

Appendices

A Lecture notes

appendix1

B Trigonometric formula tables

C Trigonometric tables