

Numerical Analysis  
Exam Minimum

Astral Projection, Yiuko

2026 年 1 月 13 日



# 目录

<b>一 数学基础知识</b>	<b>7</b>
(一) 核心概念与理论 . . . . .	7
1. 线性空间 . . . . .	7
2. 度量与赋范空间 . . . . .	8
3. 内积空间 . . . . .	11
4. 正交多项式 . . . . .	12
5. 矩阵分析回顾 . . . . .	15
6. 矩阵空间 . . . . .	16
<b>二 函数插值与重构</b>	<b>19</b>
(一) 通用理论 . . . . .	20
1. 问题模型 . . . . .	20
2. 插值空间 . . . . .	20
3. 误差分析与收敛性 . . . . .	20
(二) 具体插值方法 . . . . .	20
1. 一维多项式插值 . . . . .	20
2. 分段插值 . . . . .	25
3. Fourier 插值 . . . . .	27
<b>三 函数逼近</b>	<b>29</b>
(一) 通用理论 . . . . .	30
1. 问题模型 . . . . .	30
2. 逼近准则 . . . . .	30
3. 核心定理 . . . . .	30
(二) 具体逼近方法 . . . . .	30
1. 最优平方逼近 . . . . .	30
2. 最小二乘逼近: 最优平方逼近的离散化形式 . . . . .	33
3. 最佳一致逼近 . . . . .	34

<b>四 数值微积分</b>	<b>37</b>
(一) 数值积分 . . . . .	39
1. 通用理论 . . . . .	39
2. 具体求积方法 . . . . .	40
(二) 数值微分 . . . . .	48
1. 基础方法 . . . . .	48
2. 高精度方法 . . . . .	48
<b>五 非线性方程求根 (Nonlinear Equations)</b>	<b>51</b>
(一) 不动点迭代法 (Fixed Point Iteration) . . . . .	51
1. 基本概念 . . . . .	51
2. 收敛理论 . . . . .	51
3. 局部收敛性定理 . . . . .	52
4. 收敛阶 (Order of Convergence) . . . . .	52
5. 收敛域 (Convergence Domain) . . . . .	53
(二) 牛顿法 (Newton's Method) . . . . .	53
1. 构造 . . . . .	53
2. 收敛性分析 . . . . .	53
3. 重根情形修正 . . . . .	54
(三) 加速方法 . . . . .	54
1. 割线法 (Secant Method) . . . . .	54
2. Aitken $\Delta^2$ 加速 . . . . .	55
3. Steffensen 方法 . . . . .	55
(四) 非线性方程组 . . . . .	56
1. 基础理论与向量值微积分 . . . . .	56
2. 多元牛顿法 (Newton's Method for Systems) . . . . .	56
<b>六 常微分方程初值问题数值解法</b>	<b>59</b>
(一) 通用理论 . . . . .	60
1. 问题模型 . . . . .	60
2. 数值解法前提: Lipschitz 条件 . . . . .	60
3. 基本概念 . . . . .	60
(二) 相关定理 . . . . .	64
1. 单步法相容的充要条件 . . . . .	64
2. 整体误差估计 . . . . .	64
3. 收敛性定理 . . . . .	65
4. 稳定性定理 . . . . .	65
(三) 具体求解方法 . . . . .	66
1. 单步法 . . . . .	66

2. 多步法 . . . . .	69
<b>七 线性代数方程组数值解法</b>	<b>71</b>
(一) 线性代数方程组直接解法 (Direct Methods) . . . . .	71
1. 基本概念 . . . . .	71
2. Gauss 消去法 (Gauss Elimination) . . . . .	71
3. 矩阵三角分解 (Matrix Factorization) . . . . .	72
4. 误差分析与条件数 . . . . .	74
5. 广义逆与正则化 (Generalized Inverse & Regularization) . . . . .	76
(二) 单步定常线性迭代法 (Stationary Linear Iterative Methods) . . . . .	77
1. 迭代法基础理论 . . . . .	77
2. 经典迭代方法 . . . . .	79
(三) 非定常迭代法 (Krylov Subspace Methods) . . . . .	81
1. 变分原理与 Ritz 方法基础 . . . . .	81
2. 最速下降法 (Steepest Descent) . . . . .	83
3. 共轭梯度法 (Conjugate Gradient, CG) . . . . .	85
4. 预处理技术 (Preconditioning) . . . . .	85



# 一 数学基础知识

## (一) 核心概念与理论

### 1. 线性空间

#### (1) 定义与性质

定义 (线性空间). 设  $S$  是一个集合,  $P$  是一个数域 ( $\mathbb{R}$  或  $\mathbb{C}$ ). 定义两种映射关系:

- 向量加法:  $+ : S \times S \rightarrow S$
- 数乘:  $\cdot : P \times S \rightarrow S$

如果对任意的  $u, v, w \in S$  和  $a, b \in P$ , 满足以下八条公理, 则称  $(S, P)$  为一个线性空间 (向量空间):

1. 加法交换律:  $u + v = v + u$
2. 加法结合律:  $(u + v) + w = u + (v + w)$
3. 存在加法单位元: 存在零向量  $0 \in S$ , 使得对任意  $v \in S$ , 有  $v + 0 = v$
4. 存在加法逆元: 对任意  $v \in S$ , 存在  $-v \in S$ , 使得  $v + (-v) = 0$
5. 数乘结合律:  $a(bv) = (ab)v$
6. 数乘分配律 1:  $a(u + v) = au + av$
7. 数乘分配律 2:  $(a + b)v = av + bv$
8. 数乘单位元:  $1v = v$

则称  $(S, P)$  构成一个线性空间。

此外, 如果对于给定空间的运算法则和数域是不言自明的, 则通常简写为  $S$  是一个线性空间。如我们说  $\mathbb{R}^n$  是一个线性空间, 通常指  $(\mathbb{R}^n, \mathbb{R})$  是一个线性空间或  $(\mathbb{R}^n, \mathbb{C})$  是一个线性空间, 具体取决于数域的选择。

## (2) 线性无关与相关

待填写: (定义) 线性无关与线性相关

## (3) 基、框架与维数

待填写: (定义) 基、框架与维数

性质:

- 空间的维度是一个内蕴量, 与基的选择无关
- 多项式空间  $P_N$  中,  $\{1, x, x^2, \dots, x^N\}$  构成其一组基, 维数为  $\dim P_N = N + 1$
- 连续函数空间  $C[a, b]$  中,  $\forall N$ ,  $\{1, x, x^2, \dots, x^N\}$  是线性无关的, 但不能构成其基, 因其维数为无穷大

## 2. 度量与赋范空间

### (1) 距离空间

定义 (距离空间). 设  $M$  是一个集合,  $d : M \times M \rightarrow \mathbb{R}$  是一个映射, 如果对任意的  $x, y, z \in M$ , 满足以下三条公理, 则称  $(M, d)$  为一个距离空间:

1. 非负性与分离性:  $d(x, y) \geq 0$ , 且当且仅当  $x = y$  时,  $d(x, y) = 0$
2. 对称性:  $d(x, y) = d(y, x)$
3. 三角不等式:  $d(x, z) \leq d(x, y) + d(y, z)$

则称  $(M, d)$  构成一个距离空间。

### (2) 距离空间的完备性

待填写: (定义) 完备性

$\mathbb{R}$  是完备的, 且任意有限维赋范空间都是完备的。

#### a. 构造方法: 距离空间的完备化

设  $(M, d)$  是一个距离空间, 可以按照如下过程构造其完备化空间:

1. 构造对偶的柯西列空间

$$\tilde{M} = \{(x_n) \mid x_n \in M, (x_n) \text{ 为柯西列}\}$$

2. 在柯西列空间  $\tilde{M}$  中定义等价关系

$$\tilde{x} \sim \tilde{y} \leftrightarrow \lim_{n \rightarrow \infty} d(x_n, y_n) = 0$$

即这两个柯西列按照角标顺序，交叉放在一起，还是柯西列。

3. 构造商空间：

$$\hat{M} = \tilde{M} / \sim = \{[\tilde{x}]\}$$

式中， $[\tilde{x}]$  表示柯西列  $\tilde{x}$  的等价类，即  $[\tilde{x}]$  是一个集合，集合中的所有元素在等价关系  $\sim$  下都是等价的。

4. 在商空间  $\hat{M}$  中定义距离

$$\hat{d}([\tilde{x}], [\tilde{y}]) = \lim_{n \rightarrow \infty} d(x_n, y_n)$$

5. 则  $(\hat{M}, \hat{d})$  即为距离空间  $(M, d)$  的完备化空间。

嵌入映射：可以在原空间  $M$  与完备化空间  $\hat{M}$  之间定义一个单射  $i$ ：

$$i : M \rightarrow \hat{M}, \quad i(x) = [(x, x, x, \dots)]$$

该映射将原空间中的每个点  $x$  映射为完备化空间中由常值序列  $(x, x, x, \dots)$  所构成的等价类，且映射前后任意两元素的距离不变

### (3) 赋范空间与 Banach 空间

**待填写：(定义) 赋范空间**

完备的赋范空间称为 Banach 空间，或者 B 空间。

### (4) 等价范数

定义 (范数的等价性). 设  $V$  是一个线性空间， $\|\cdot\|_1$  和  $\|\cdot\|_2$  是  $V$  上的两个范数，如果存在正常数  $c$  和  $C$ ，使得对任意  $v \in V$ ，都有

$$c\|v\|_1 \leq \|v\|_2 \leq C\|v\|_1$$

则称这两个范数是等价的。

若存在正数  $C$ ，使得对任意  $v \in V$ ，都有

$$\|v\|_2 \leq C\|v\|_1$$

则称范数  $\|\cdot\|_1$  强于  $\|\cdot\|_2$ 。

性质：

- 在有限维线性空间上，任意两个范数都是等价的
- 在无限维线性空间上，范数不一定是等价的
- 若一个点列在较强的范数下是 Cauchy 列，则在较弱的范数下也是 Cauchy 列；反之不必然。

## (5) 常用的范数

### a. $\mathbb{R}^n$ 上的范数

记  $x = (x_1, x_2, \dots, x_n)^T \in \mathbb{R}^n$ , 则常用的范数有:

- 无穷范数: 所有元素的最大值

$$\|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$$

- 1-范数: 所有元素的绝对值之和

$$\|x\|_1 = \sum_{i=1}^n |x_i|$$

- 2-范数: 欧几里得范数, 即所有元素的平方和的平方根

$$\|x\|_2 = \left( \sum_{i=1}^n |x_i|^2 \right)^{\frac{1}{2}}$$

### b. $C[a, b]$ (有界闭区间上连续函数空间) 上的范数

- 无穷范数: 函数在区间上的最大绝对值

$$\|f\|_\infty = \max_{a \leq x \leq b} |f(x)|$$

- 1-范数: 函数在区间上的绝对值积分

$$\|f\|_1 = \int_a^b |f(x)| dx$$

- 2-范数: 函数在区间上的平方积分的平方根

$$\|f\|_2 = \left( \int_a^b |f(x)|^2 dx \right)^{\frac{1}{2}}$$

$C[a, b]$  上三个范数的性质:

- 任意两个范数不等价
- 无穷范数强于 2 范数, 2 范数强于 1 范数
- 只有无穷范数对应的赋范空间是完备的
- 1-范数对应的完备化空间为  $L^1(a, b)$ , 2-范数对应的完备化空间为  $L^2(a, b)$ <sup>1</sup>

<sup>1</sup>  $L^1(a, b)$  为  $(a, b)$  上的可积函数空间,  $L^2(a, b)$  为  $(a, b)$  上的平方可积函数空间。

### 3. 内积空间

#### (1) 内积定义与性质

**定义 (内积).** 设  $(S, P)$  是一个线性空间, 如果对任意的  $u, v, w \in S$  和  $a, b \in P$ , 存在一个映射  $S \times S \rightarrow P$ , 满足

1. 共轭对称性:  $\langle u, v \rangle = \overline{\langle v, u \rangle}$
2. 线性性:  $\langle au + bv, w \rangle = a\langle u, w \rangle + b\langle v, w \rangle$
3. 正定性:  $\langle v, v \rangle \geq 0$ , 且当且仅当  $v = 0$  时,  $\langle v, v \rangle = 0$

则称该映射为内积,  $(S, P)$  构成一个内积空间。

若  $\langle x, y \rangle = 0$ , 则称  $x$  与  $y$  正交。

几个常用空间上的内积:

- $\mathbb{R}^n$  或  $\mathbb{C}^n$  上的内积

$$\langle x, y \rangle = \sum_{i=1}^n x_i y_i$$

- $C[a, b]$  (有界闭区间上连续函数空间) 上的内积

$$\langle f, g \rangle = \int_a^b f(x) \overline{g(x)} dx$$

- $C[a, b]$  上的带权内积:

$$(f, g) = \int_a^b \rho(x) f(x) \overline{g(x)} dx$$

其中, 权函数  $\rho(x)$  需要满足条件:

- $\rho(x) \in C[a, b]$
- $\rho(x)$  几乎处处为正
- $\int_a^b \rho(x) dx < +\infty$
- $\forall q(x) \in P_n$ ,  $\int_a^b \rho(x) |q(x)| dx < \infty$

带权内积所研究的空间称为加权内积空间:

$$L_\rho^2(a, b) = \{f(x) \mid \int_a^b \rho(x) |f(x)|^2 dx < +\infty\}$$

常用的权函数有:

$$\begin{aligned} \rho(x) &= 1, \quad [a, b] = [-1, 1] \\ \rho(x) &= \frac{1}{1 - x^2}, \quad [a, b] = [-1, 1] \end{aligned}$$

## (2) 正交性与 Schmidt 正交化

待填写: (定义) 正交性

待填写: (方法) 用 Grammer 矩阵判断内积空间中向量组的线性无关性

待填写: (方法) Schmidt 正交化过程: 从一个线性无关向量组构造一个正交向量组: 让每个向量减去与已有空间垂直的分量

用 Schmidt 正交化过程得到的正交向量组具有以下性质:

$$\Phi_{k-1} \subset \Phi_k$$

$$y_k \perp \Phi_{k-1}$$

## (3) 由内积诱导的范数

定义 (诱导范数). 设  $(S, P)$  是一个内积空间, 则可以定义范数  $\|\cdot\|$  如下:

$$\|v\| = \sqrt{\langle v, v \rangle}$$

则称该范数为由内积诱导的范数。

- 任何内积均能诱导对应的范数
- 当且仅当范数满足平行四边形法则时

$$\|f + g\|^2 + \|f - g\|^2 = 2\|f\|^2 + 2\|g\|^2$$

范数可以诱导内积:

$$(x, y) = \frac{1}{4}(\|x + y\|^2 - \|x - y\|^2 + i\|x + iy\|^2 - i\|x - iy\|^2)$$

## 4. 正交多项式

定义 (正交多项式). 设  $\{\phi_n(x)\}$  是定义在区间  $[a, b]$  上的一组多项式, 且每个多项式的次数为  $n$ , 如果对任意  $m \neq n$ , 都有

$$\int_a^b \rho(x) \phi_m(x) \phi_n(x) dx = 0$$

则称  $\{\phi_n(x)\}$  为区间  $[a, b]$  上关于权函数  $\rho(x)$  的正交多项式。

正交多项式的性质:

- $\deg \phi_i = i$
- $(\phi_i, \phi_j) = 0, \quad \forall i \neq j$
- $\phi_n$  为实系数多项式
- $\phi_n$  在开区间  $(a, b)$  内恰有  $n$  个实单根

待填写: (证明)  $\phi_n$  在开区间  $(a, b)$  内恰有  $n$  个实单根的证明。证法: 分别证明实根、单根、全在  $(a, b)$  内。3 个命题均可用反证法。

(1)  $\rho = 1$ : Legendre 多项式

产生方法:

- 权函数  $\rho(x) = 1$
- 区间  $[-1, 1]$

表达式:

$$P_n(x) = \frac{1}{2^n n!} \frac{d^n}{dx^n} [(x^2 - 1)^n]$$

性质:

- 首项系数:

$$k_n = \frac{(2n)!}{2^n (n!)^2}$$

- 正交归一化:

$$\int_{-1}^1 P_n(x) P_m(x) dx = \frac{2}{2n+1} \delta_{mn}$$

- 三项递推关系:

$$(n+1)P_{n+1}(x) = (2n+1)xP_n(x) - nP_{n-1}(x)$$

- 奇偶性:

$$P_n(-x) = (-1)^n P_n(x)$$

- 导数关系:

$$\frac{d}{dx} P_n(x) = \frac{n}{x^2 - 1} [xP_n(x) - P_{n-1}(x)]$$

- 前五项:

$$P_0(x) = 1$$

$$P_1(x) = x$$

$$P_2(x) = \frac{1}{2}(3x^2 - 1)$$

$$P_3(x) = \frac{1}{2}(5x^3 - 3x)$$

$$P_4(x) = \frac{1}{8}(35x^4 - 30x^2 + 3)$$

- 零平方误差最小:

**定理.** 在所有首项为 1 的  $n$  次多项式中, Legendre 多项式  $\tilde{P}_n(x)$  在  $[-1, 1]$  上与零的平方误差最小。

(2)  $\rho(x) = \frac{1}{\sqrt{1-x^2}}$ : Chebyshev 多项式

产生方法:

- 权函数  $\rho(x) = \frac{1}{\sqrt{1-x^2}}$
- 区间  $[-1, 1]$

表达式:

$$T_n(x) = \cos(n \arccos x)$$

性质:

- 首项系数:  $2^{n-1}$
- 正交归一化:

$$\int_{-1}^1 \frac{T_n(x)T_m(x)}{\sqrt{1-x^2}} dx = \begin{cases} \pi, & n = m = 0 \\ \frac{\pi}{2}, & n = m \neq 0 \\ 0, & n \neq m \end{cases}$$

- 三项递推关系:

$$T_{n+1} = 2xT_n - T_{n-1}$$

- 奇偶性:

$$T_n(-x) = (-1)^n T_n(x)$$

- 前五项:

$$\begin{aligned} T_0(x) &= 1 \\ T_1(x) &= x \\ T_2(x) &= 2x^2 - 1 \\ T_3(x) &= 4x^3 - 3x \\ T_4(x) &= 8x^4 - 8x^2 + 1 \end{aligned}$$

- 零点:

$$x_k = \cos\left(\frac{2k-1}{2n}\pi\right), \quad k = 1, 2, \dots, n$$

- 极值点:

$$x_k = \cos\left(\frac{k\pi}{n}\right), \quad k = 0, 1, \dots, n$$

- 简单表达式: 当  $|x| \geq 1$  时,

$$T_n(x) = \frac{1}{2} \left[ \left( x + \sqrt{x^2 - 1} \right)^n + \left( x - \sqrt{x^2 - 1} \right)^n \right]$$

## 5. 矩阵分析回顾

### (1) 矩阵代数与初等变换 (Elementary Transformations)

初等矩阵：由单位矩阵  $I$  经过一次初等变换得到的矩阵。

1. 交换矩阵 (Permutation Matrix)  $P_{ij}$ : 交换  $I$  的第  $i$  行和第  $j$  行。

- 作用：左乘交换矩阵的行，右乘交换矩阵的列。
- 性质： $P_{ij}^{-1} = P_{ij}^T = P_{ij}$  (对称且正交)。
- 应用：Gauss 消去法中的选主元 ( $PA = LU$ )。

2. 倍乘矩阵 (Scaling Matrix)  $D_i(k)$ : 将  $I$  的第  $i$  行乘以非零常数  $k$ 。

- 性质： $D_i(k)$  是对角阵。 $D_i(k)^{-1} = D_i(1/k)$ 。

3. 倍加矩阵 (Elimination/Shear Matrix)  $E_{ij}(k)$  ( $i \neq j$ ): 将  $I$  的第  $j$  行的  $k$  倍加到第  $i$  行。

$$E_{ij}(k) = I + ke_i e_j^T$$

- 作用：左乘  $A$  将第  $j$  行的  $k$  倍加到第  $i$  行 (消元操作)。
- 性质：1. 逆矩阵形式简单： $(E_{ij}(k))^{-1} = E_{ij}(-k)$ 。2. Gauss 消去法中的消元矩阵  $L_k$  就是一系列倍加矩阵的乘积。

### (2) 矩阵特征值理论

- 特征分解： $Ax = \lambda x$ 。
- SVD 分解： $A = U\Sigma V^H$ 。
- Schur 分解： $U^H A U = T$ 。

### (3) 特殊矩阵类

- Hermite 矩阵 (实对称矩阵)： $A^H = A$ 。
  - 特征值全为实数。
  - 存在酉矩阵  $U$  使得  $U^H A U = \Lambda$  (可酉对角化)。
  - Rayleigh 商： $\lambda_{\min} \leq \frac{x^H A x}{x^H x} \leq \lambda_{\max}$ 。
- 对称正定矩阵 (Symmetric Positive Definite, SPD)：

定义 (SPD 定义)。实对称矩阵  $A$  称为正定的，如果对任意非零向量  $x \in \mathbb{R}^n, x \neq 0$ ，都有：

$$x^T A x > 0$$

这定义了一个能量范数 (Energy Norm)： $\|x\|_A = \sqrt{x^T A x}$ 。

性质：

1. 特征值：所有特征值均严格大于 0 ( $\lambda_i > 0$ )。
2. 行列式： $\det(A) = \prod \lambda_i > 0$ , 故  $A$  必可逆。
3. 主子式：所有顺序主子式均大于 0 (Sylvester 准则)。
4. 对角元： $a_{ii} > 0$ , 且  $\max_{i,j} |a_{ij}| = \max_i a_{ii}$  (最大元素必在对角线上)。
5. 逆矩阵：若  $A$  是 SPD, 则  $A^{-1}$  也是 SPD。
6. Cholesky 分解： $A$  存在唯一的分解  $A = LL^T$ , 其中  $L$  为对角元为正的下三角阵。
- 酉矩阵(正交矩阵)： $U^H U = I$ 。保持向量 2-范数不变 ( $\|Ux\|_2 = \|x\|_2$ )。

## 6. 矩阵空间

待填写：(性质) 矩阵空间的基本性质：线性空间、乘法运算、代数性质

### (1) 矩阵范数

定义(矩阵范数). 矩阵空间  $\mathbb{C}^{n \times n}$  上的范数  $\|\cdot\|$  称为矩阵范数, 如果对任意的  $A, B \in \mathbb{C}^{n \times n}$  和  $a \in \mathbb{C}$ , 满足以下性质：

1. 非负性与分离性： $\|A\| \geq 0$ , 且当且仅当  $A = 0$  时,  $\|A\| = 0$
2. 齐次性： $\|aA\| = |a|\|A\|$
3. 三角不等式： $\|A + B\| \leq \|A\| + \|B\|$
4. 次乘性： $\|AB\| \leq \|A\| \cdot \|B\|$

Note: 矩阵范数是定义在矩阵代数而非矩阵空间上的, 必须与矩阵乘法相容。

定义(矩阵范数与向量范数的相容性). 设  $\|\cdot\|_v$  是向量空间  $\mathbb{C}^n$  上的一个范数,  $\|\cdot\|_m$  是矩阵空间  $\mathbb{C}^{n \times n}$  上的一个范数, 如果对任意的  $A \in \mathbb{C}^{n \times n}$  和  $x \in \mathbb{C}^n$ , 都有

$$\|Ax\|_v \leq \|A\|_m \cdot \|x\|_v$$

则称矩阵范数  $\|\cdot\|_m$  与向量范数  $\|\cdot\|_v$  是相容的。

矩阵范数的两种常见构造方法：

- 直接构造: Frobenius 范数

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

Frobenius 范数的性质：

- Frobenius 范数与向量 2-范数相容
- $\|I\|_F = \sqrt{n}$
- 向量范数诱导: 算子范数

**定义 (算子范数).** 设  $\|\cdot\|_v$  是向量空间  $\mathbb{C}^n$  上的一个范数, 则可以定义矩阵空间  $\mathbb{C}^{n \times n}$  上的算子范数  $\|\cdot\|_m$  如下:

$$\|A\|_m = \max_{x \neq 0} \frac{\|Ax\|_v}{\|x\|_v} = \max_{\|x\|_v=1} \|Ax\|_v$$

则称该范数为由向量范数  $\|\cdot\|_v$  诱导的算子范数。

常用的几个算子范数:

- 无穷范数: 行和最大值

$$\|A\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|$$

- 1-范数: 列和最大值

$$\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}|$$

- 2-范数 (谱范数):  $A$  的最大奇异值, 即  $A^H A$  的最大特征值的平方根

$$\|A\|_2 = \sqrt{\lambda_{\max}(A^H A)}$$

## (2) 谱半径

**定义 (谱半径).** 谱半径定义为矩阵所有特征值模的最大值, 即

$$\rho(A) = \max_{1 \leq i \leq n} |\lambda_i|$$

谱半径和矩阵范数的关系:

- 矩阵范数下界:

**定理.** 对任意  $A \in \mathbb{C}^{n \times n}$ , 有

$$\rho(A) \leq \|A\|$$

- 无穷接近范数的存在性:

**定理.** 对任意  $A \in \mathbb{C}^{n \times n}$ , 存在一个矩阵范数  $\|\cdot\|$ , 使得

$$\rho(A) \leq \|A\| \leq \rho(A) + \varepsilon$$

其中,  $\varepsilon$  为任意给定的正常数。

### (3) 可逆矩阵相关定理

**定理 (扰动引理 I).** 给定  $B \in \mathbb{C}^{n \times n}$ 。设  $\|B\| < 1$ , 则  $I + B$  可逆, 且

$$\|(I + B)^{-1}\| \leq \frac{1}{1 - \|B\|}$$

**定理 (扰动引理 II).** 设  $A, C \in \mathbb{C}^{n \times n}$ , 且  $A$  可逆。若

$$\|C - A\| < \frac{1}{\|A^{-1}\|}$$

则  $C$  也可逆, 且

$$\|C^{-1}\| \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \cdot \|C - A\|}$$

**定理 (扰动定理 II).** 设  $A, \delta A \in \mathbb{C}^{n \times n}$ , 且  $A$  可逆。若  $\|A^{-1}\delta A\| < 1$ , 则  $A + \delta A$  也可逆, 且

$$\|(A + \delta A)^{-1}\| \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \cdot \|\delta A\|}$$

## 二 函数插值与重构

总结：基本方法是利用插值基函数构造插值多项式，从而实现对函数的近似与重构。

问题 1：求解插值函数

- 整个区间上的连续插值：Lagrange 插值、Newton 插值、Hermite 插值

– Lagrange 插值基函数：

$$L_\alpha(x) = \prod_{\substack{\beta \in I \\ \beta \neq \alpha}} \frac{x - x_\beta}{x_\alpha - x_\beta}, \quad \alpha \in I$$

– Newton 插值和 Hermite 插值：构造均差表，列表计算。

\* 均差递推关系：

$$\begin{aligned} f[x_i] &= f(x_i) \\ f[x_i, x_{i+1}, \dots, x_{i+k}] &= \frac{f[x_{i+1}, \dots, x_{i+k}] - f[x_i, \dots, x_{i+k-1}]}{x_{i+k} - x_i} \end{aligned}$$

\* 均差构造插值多项式：各项为  $f[x_0, x_1, \dots, x_k] \cdot (x - x_0)(x - x_1) \dots (x - x_{k-1})$

- 分段插值：分片线性插值、分片三次 Hermite 插值

– 分片线性插值基函数：

$$L_\alpha(x) = \begin{cases} \frac{x - x_{\alpha-1}}{x_\alpha - x_{\alpha-1}}, & x \in [x_{\alpha-1}, x_\alpha] \\ \frac{x_{\alpha+1} - x}{x_{\alpha+1} - x_\alpha}, & x \in [x_\alpha, x_{\alpha+1}] \\ 0, & \text{else} \end{cases}$$
$$\phi(x) = \sum_{\alpha=0}^n f_\alpha L_\alpha(x)$$

- 分片三次 Hermite 插值基函数：

$$\begin{aligned}\alpha_k &= \left(1 + 2\frac{x - x_k}{x_{k+1} - x_k}\right) \left(\frac{x_{k+1} - x}{x_{k+1} - x_k}\right)^2 \\ \beta_k &= (x - x_k) \left(\frac{x_{k+1} - x}{x_{k+1} - x_k}\right)^2 \\ \alpha_{k+1} &= \left(1 + 2\frac{x_{k+1} - x}{x_{k+1} - x_k}\right) \left(\frac{x - x_k}{x_{k+1} - x_k}\right)^2 \\ \beta_{k+1} &= -(x_{k+1} - x) \left(\frac{x - x_k}{x_{k+1} - x_k}\right)^2 \\ \phi(x) &= \sum_{k=0}^n [f_k \alpha_k + f'_k \beta_k], \quad x \in [x_k, x_{k+1}]\end{aligned}$$

## (一) 通用理论

### 1. 问题模型

待填写：(数学描述) 采样泛函视角下的插值问题数学描述

### 2. 插值空间

常用的插值空间：多项式函数空间、样条函数空间、三角多项式函数空间

### 3. 误差分析与收敛性

## (二) 具体插值方法

### 1. 一维多项式插值

问题：给定插值数据（采样数据） $(x_\alpha, f_\alpha)$ ,  $\alpha \in I$ , 确定多项式  $P(x) \in P_n$ ,  $n = |I| - 1$ , 满足插值条件

$$x_\alpha(P) = P(x_\alpha) = f_\alpha, \quad \alpha \in I$$

定理 (多项式插值基本定理). 给定  $n + 1$  个插值条件

$$(x_\alpha, f_\alpha), \quad \alpha \in I, \quad x_\alpha \neq x_\beta \text{ for } \alpha \neq \beta$$

则存在唯一的插值多项式  $P \in P_n$  满足插值条件。

Note: 若  $x_\alpha$  取之于复平面, 上述定理依然成立; 且上述定理与采样节点的排序无关。

## (1) Lagrange 插值

### a. 基函数构造

定义 (Lagrange 插值基函数). 定义

$$L_\alpha(x) = L_{\alpha;I}(x) = \prod_{\substack{\beta \in I \\ \beta \neq \alpha}} \frac{x - x_\beta}{x_\alpha - x_\beta}, \quad \alpha \in I$$

称为插值基函数。

若给定 3 个插值条件  $(x_0, f_0), (x_1, f_1), (x_2, f_2)$ , 则对应的插值基函数为

$$\begin{aligned} L_0(x) &= \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} \\ L_1(x) &= \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)} \\ L_2(x) &= \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)} \end{aligned}$$

插值基函数天然满足性质:

$$x_\beta(L_\alpha) = L_\alpha(x_\beta) = \delta_{\alpha\beta}$$

### b. 插值公式

在计算出插值基函数的基础上, 插值多项式可写为:

$$P(x) = \sum_{\alpha \in I} f_\alpha L_\alpha(x)$$

### c. 余项

### d. 均差定义

定义 (均差的递推公式). 设  $f(x)$  在区间  $[a, b]$  上有定义, 且给定插值节点  $x_0, x_1, \dots, x_n$ , 则定义如下均差:

$$\begin{aligned} f[x_i] &= f(x_i) \\ f[x_i, x_{i+1}, \dots, x_{i+k}] &= \frac{f[x_{i+1}, \dots, x_{i+k}] - f[x_i, \dots, x_{i+k-1}]}{x_{i+k} - x_i} \end{aligned}$$

其中,  $i = 0, 1, \dots, n - k$ ,  $k = 1, 2, \dots, n$ .

均差的性质:

- $f_{i_0 i_1 \dots i_k}$  与节点  $x_{i_0}, x_{i_1}, \dots, x_{i_k}$  的顺序无关
- 设  $f$  是  $N$  次多项式, 若  $k > N$ , 则对任意节点  $x_{i_0}, x_{i_1}, \dots, x_{i_k}$ , 都有

$$f_{i_0 i_1 \dots i_k} = 0$$

## e. 插值公式

$$P_{i_0 i_1 \dots i_k}(x) = f_{i_0} + f_{i_0 i_1}(x - x_{i_0}) + f_{i_0 i_1 i_2}(x - x_{i_0})(x - x_{i_1}) + \dots + f_{i_0 i_1 \dots i_k}(x - x_{i_0})(x - x_{i_1}) \cdots (x - x_{i_{k-1}})$$

## f. Newton 插值多项式的列表计算

以给定 4 个节点时  $x_0, x_1, x_2, x_3$  的插值问题为例。可以按照如下表格从左向右逐列填写计算均差：

$x_i$	0 阶均差	1 阶均差	2 阶均差	3 阶均差
$x_0$	$f(x_0)$	$f[x_0, x_1] = \frac{f(x_1) - f(x_0)}{x_1 - x_0}$		
$x_1$	$f(x_1)$		$f[x_0, x_1, x_2] = \frac{f[x_1, x_2] - f[x_0, x_1]}{x_2 - x_0}$	
$x_2$	$f(x_2)$	$f[x_1, x_2] = \frac{f(x_2) - f(x_1)}{x_2 - x_1}$	$f[x_1, x_2, x_3] = \frac{f[x_2, x_3] - f[x_1, x_2]}{x_3 - x_1}$	$f[x_0, x_1, x_2, x_3] = \dots$
$x_3$	$f(x_3)$	$f[x_2, x_3] = \frac{f(x_3) - f(x_2)}{x_3 - x_2}$		

之后用插值表最上方一行的均差值逐个组装 Newton 插值多项式。

## (2) Hermite 插值

## a. Hermite 插值问题

给定  $\xi_i, f_i^{(k)}, i = 0, 1, \dots, m, k = 0, 1, \dots, n_i - 1$ , 其中  $\xi_i$  两两不同, 且

$$\xi_0 < \xi_1 < \dots < \xi_m$$

希望确定一个次数为  $n$  的多项式函数

$$P_n(x), \quad n = \sum_{i=0}^m n_i - 1$$

满足插值条件

$$P^{(k)}(\xi_i) = f_i^{(k)}, \quad i = 0, 1, \dots, m \quad k = 0, 1, \dots, n_i - 1$$

### b. 拓展均差

**定义 (拓展均差).** 设  $f \in C^n(I(x_0, x_1, \dots, x_n))$ , 定义

$$f[x_0, x_1, x_n] = \int_0^{t_0} dt_1 \int_0^{t_1} dt_2 \cdots \int_0^{t_{n-1}} dt_n$$

$$f^{(n)}(t_n[x_n - x_{n-1}] + t_{n-1}[x_{n-1} - x_{n-2}] + \dots + t_1[x_1 - x_0] + t_0 x_0)$$

式中,  $n \geq 1$  且  $t_0 = 1$ 。

Note: 这一积分实际上表示了一个单位标准  $n$ -维标准型上的积分, 或者说积分区域始终是一个插值节点构造的凸组合。这隐含了一个要求是  $1 = t_0 \geq t_1 \geq t_2 \geq \dots \geq t_n \geq 0$ 。

拓展均差的性质:

- 若  $x_i$  两两不一, 则拓展均差等价于普通均差

$$f[x_0, x_1, x_n] = f_{x_0, x_1, \dots, x_n}$$

且具有相同的递推关系:

$$f[x_0, x_1, \dots, x_n] = \frac{f[x_0, x_1, \dots, x_{n-2}, x_n] - f[x_0, x_1, \dots, x_{n-1}]}{x_n - x_0}$$

- 极限性质: 若  $f$  足够光滑, 则

$$\lim_{\epsilon_i \rightarrow 0} f[x_0 + \epsilon_0, x_1 + \epsilon_1, \dots, x_n + \epsilon_n] = f[x_0, x_1, \dots, x_n]$$

- 导数与重节点: 可从极限性质导出

$$\frac{d}{dx} f[x_0, x_1, \dots, x_n, x] = f[x_0, x_1, \dots, x_n, x, x]$$

- 介值定理: 若  $f \in C^n[a, b]$ ,  $x_0, x_1, \dots, x_n \in [a, b]$ , 则存在  $\xi \in I(x_0, x_1, \dots, x_n)$ , 使得

$$f[x_0, x_1, \dots, x_n] = \frac{f^{(n)}(\xi)}{n!}$$

特别地,

$$\underbrace{f[x, x, \dots, x]}_{n+1 \uparrow} = \frac{f^{(n)}(x)}{n!}$$

### c. Hermite 插值多项式

Hermite 插值多项式可表示为

$$P(x) = f[x_0] + f[x_0, x_1](x - x_0) + \dots + f[x_0, x_1, \dots, x_n](x - x_0)(x - x_1)\dots(x - x_{n-1})$$

其中,  $x_0, \dots, x_n$  为下面序列的任意置换:

$$\underbrace{\xi_0, \xi_0, \dots, \xi_0}_{n_0 \uparrow}, \underbrace{\xi_1, \xi_1, \dots, \xi_1}_{n_1 \uparrow}, \dots, \underbrace{\xi_m, \xi_m, \dots, \xi_m}_{n_m \uparrow}$$

#### d. 列表法求 Hermite 插值多项式

假设给定 2 个节点  $\xi_0, \xi_1$ , 对应的插值条件分别为  $f_0, f'_0, f''_0, f_1$ , 则可按下表计算均差:

Hermite 插值均差表 (节点序列:  $\xi_0, \xi_0, \xi_0, \xi_1$ ):

节点	0 阶均差	1 阶均差	2 阶均差	3 阶均差
$\xi_0$	$f[\xi_0] = f_0$			
$\xi_0$	$f[\xi_0] = f_0$	$f[\xi_0, \xi_0] = f'_0$		
$\xi_0$	$f[\xi_0] = f_0$	$f[\xi_0, \xi_0] = f'_0$	$f[\xi_0, \xi_0, \xi_0] = \frac{f''_0}{2}$	
$\xi_1$	$f[\xi_1] = f_1$	$f[\xi_0, \xi_1] = \frac{f_1 - f_0}{\xi_1 - \xi_0}$	$f[\xi_0, \xi_0, \xi_1] = \frac{f[\xi_0, \xi_1] - f[\xi_0, \xi_0]}{\xi_1 - \xi_0}$	$f[\xi_0, \xi_0, \xi_0, \xi_1] = \frac{f[\xi_0, \xi_0, \xi_1] - f[\xi_0, \xi_0, \xi_0]}{\xi_1 - \xi_0}$

对应的插值多项式为:

$$\begin{aligned} P(x) &= f[\xi_0] + f[\xi_0, \xi_0](x - \xi_0) + f[\xi_0, \xi_0, \xi_0](x - \xi_0)^2 + f[\xi_0, \xi_0, \xi_0, \xi_1](x - \xi_0)^3 \\ &= f_0 + f'_0(x - \xi_0) + \frac{f''_0}{2}(x - \xi_0)^2 \\ &\quad + \frac{2(f_1 - f_0 - f'_0(\xi_1 - \xi_0)) - f''_0(\xi_1 - \xi_0)^2}{2(\xi_1 - \xi_0)^3}(x - \xi_0)^3 \end{aligned}$$

### (3) Lagrange 插值和 Hermite 插值的收敛分析

#### a. 插值余项

若  $f$  在区间  $[a, b]$  上具有  $n+1$  阶连续导数, 则对任意  $x \in [a, b]$ , 存在  $\xi \in I(x, x_0, x_1, \dots, x_n)$ , 使得

$$R(x) = f(x) - P(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \prod_{i=0}^n (x - x_i) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \omega_{01\dots n}(x)$$

此外, 使用  $n+1$  个节点  $x_0, x_1, \dots, x_n$  进行 Hermite 插值时, 总有

$$R(x) = f(x) - P(x) = f[x_0, x_1, \dots, x_n, x](x - x_0)(x - x_1) \cdots (x - x_n)$$

#### b. 收敛性

收敛性的定义: 当给定插值点的最大间距  $h \rightarrow 0$  时, 插值余项  $R(x) \rightarrow 0$ , 则称插值多项式序列在区间  $[a, b]$  上收敛于函数  $f(x)$ 。

定义 (多项式插值收敛定义). 设  $f \in C^\infty[a, b]$ , 且存在正常数  $M > 0$ , 使得对任意的  $n \geq 0$ , 都有

$$\max_{x \in [a, b]} |f^{(n)}(x)| \leq M$$

则对任意在  $[a, b]$  上的插值节点序列  $\{x_i^{(n)}\}_{i=0}^n$ , 对应的插值多项式序列  $\{P_n(x)\}$  在  $[a, b]$  上均匀收敛于  $f(x)$ 。

收敛的充分条件：

**定理.** 记  $\delta = |I(x_0, x_1, x_2, \dots, x_n)|$ ,  $\tilde{x}$  为  $I$  的中心。若  $f$  在  $B(\tilde{x}, w\delta)$  上复解析, 则对任意  $\bar{x} \in I$ , 插值法收敛。

## 2. 分段插值

### (1) 分段线性插值

待填写：() 分片线性插值问题描述

#### a. 分片线性插值的插值基函数

**定义** (分片线性插值基函数). 设给定插值节点  $x_0 < x_1 < \dots < x_n$ , 则定义分片线性插值基函数为

$$l_i(x) = \begin{cases} \frac{x-x_{i-1}}{x_i-x_{i-1}}, & x \in [x_{i-1}, x_i] \\ \frac{x_{i+1}-x}{x_{i+1}-x_i}, & x \in [x_i, x_{i+1}] \\ 0, & \text{otherwise} \end{cases}$$

其中,  $i = 0, 1, \dots, n$ , 且约定  $x_{-1} = x_0$ ,  $x_{n+1} = x_n$ 。

则线性插值基函数为

$$\phi(x) = \sum_{k=1}^n f_k l_k(x)$$

#### b. 分片线性插值的收敛性定理

定义

$$h = \max_{1 \leq i \leq n} (x_i - x_{i-1})$$

- 若  $f \in C[a, b]$ , 则  $\lim_{h \rightarrow 0} \|f - \phi\|_\infty \rightarrow 0$
- 若  $f \in C^1[a, b]$ , 则  $\|f - \phi\|_\infty \leq \frac{h}{2} \|f'\|_\infty$
- 若  $f \in C^2[a, b]$ , 则  $\|f - \phi\|_\infty \leq \frac{h^2}{8} \|f''\|_\infty$

### (2) 分段三次 Hermite 插值

待填写：() 分片三次 Hermite 插值的数学描述

### a. 插值基函数

分段三次 Hermite 插值的基函数满足如下条件:

$$\begin{aligned}\alpha_k(x_i) &= \delta_{ik}, \quad \alpha'_k(x_i) = 0 \\ \beta_k(x_i) &= 0, \quad \beta'_k(x_i) = \delta_{ik}\end{aligned}$$

在单元  $[x_k, x_{k+1}]$  上, 三次 Hermite 插值基函数为

$$\begin{aligned}\alpha_k &= \left(1 + 2\frac{x - x_k}{x_{k+1} - x_k}\right) \left(\frac{x_{k+1} - x}{x_{k+1} - x_k}\right)^2 \\ \beta_k &= (x - x_k) \left(\frac{x_{k+1} - x}{x_{k+1} - x_k}\right)^2 \\ \alpha_{k+1} &= \left(1 + 2\frac{x_{k+1} - x}{x_{k+1} - x_k}\right) \left(\frac{x - x_k}{x_{k+1} - x_k}\right)^2 \\ \beta_{k+1} &= -(x_{k+1} - x) \left(\frac{x - x_k}{x_{k+1} - x_k}\right)^2\end{aligned}$$

- $\alpha_k$  满足单位分解性:  $\alpha_k + \alpha_{k+1} = 1$
- $\alpha_k(x_i) = \delta_{ik}$ ,  $\alpha'_k(x_i) = 0$
- $\beta_k(x_i) = 0$ ,  $\beta'_k(x_i) = \delta_{ik}$

### b. 三次 Hermite 插值多项式

$$\phi = \sum_{k=0}^n [f_k \alpha_k(x) + f'_k \beta_k(x)]$$

### c. 收敛性定理

定义 (分段三次 Hermite 插值收敛性定理). 设  $f \in C^1[a, b]$ , 则分段三次 Hermite 插值多项式  $\phi$  满足

$$\|f - \phi\|_\infty \leq ch \|f'\|_\infty$$

若  $f$  有更好的光滑性, 则:

- 若  $f \in C^2[a, b]$ , 则  $\|f - \phi\|_\infty \leq ch^2 \|f''\|_\infty$
- 若  $f \in C^3[a, b]$ , 则  $\|f - \phi\|_\infty \leq ch^3 \|f'''\|_\infty$
- 若  $f \in C^4[a, b]$ , 则  $\|f - \phi\|_\infty \leq \frac{1}{384} h^4 \|f^{(4)}\|_\infty$

此外, 对于不高于三次的多项式, 分段三次 Hermite 插值是精确的。

### 3. Fourier 插值

#### (1) 离散傅里叶变换

待填写: (定义) 离散傅里叶变换式、变换式系数表达式

如果  $f$  的光滑性满足  $f \in C_{per}^M$ , 则有

$$\begin{aligned} a_n &= O(n^{-M}) \\ b_n &= O(n^{-M}) \end{aligned}$$

且

$$\|f(x) - \left[ \frac{a_0}{2} + \sum_{n=1}^N a_n \cos nx + \sum_{n=1}^N b_n \sin nx \right] \|_\infty = O(N^{-M})$$

#### (2) 三角多项式插值空间

三角多项式插值空间为:

$$\begin{aligned} \Phi_{2M+1} &:= \left\{ \frac{A_0}{2} + \sum_{n=1}^M (A_n \cos nx + B_n \sin nx) \right\} \\ \Phi_{2M} &:= \left\{ \frac{A_0}{2} + \sum_{n=1}^{M-1} (A_n \cos nx + B_n \sin nx) + A_M \cos Mx \right\} \end{aligned}$$

#### (3) 三角多项式插值与一般多项式插值

待填写: (问题) 三角多项式插值问题的数学描述

待填写: (问题) 辅助插值问题: 找相多项式

待填写: (理论) 两个插值问题之间的联系: 欧拉公式

#### (4) 插值定理与三角插值多项式

定理 (三角多项式插值定理). 设给定插值节点

$$x_k = \frac{2k\pi}{N}, \quad k = 0, 1, \dots, N-1$$

则对任意插值数据  $f_k$ , 存在唯一的三角多项式

$$P(x) = \frac{A_0}{2} + \sum_{n=1}^M (A_n \cos nx + B_n \sin nx)$$

(当  $N$  为奇数时,  $M = \frac{N-1}{2}$ ; 当  $N$  为偶数时,  $M = \frac{N}{2}$ ) 满足插值条件

$$P(x_k) = f_k, \quad k = 0, 1, \dots, N-1$$

同样，存在唯一的相多项式

$$Q(x) = \sum_{n=0}^{N-1} \beta_n e^{inx}$$

满足插值条件

$$Q(x_k) = f_k, \quad k = 0, 1, \dots, N - 1$$

- 三角插值多项式：

$$A_j = \frac{2}{N} \sum_{k=0}^{N-1} f_k \cos\left(\frac{2\pi j k}{N}\right), \quad j = 0, 1, \dots, M$$

$$B_j = \frac{2}{N} \sum_{k=0}^{N-1} f_k \sin\left(\frac{2\pi j k}{N}\right), \quad j = 1, 2, \dots, M$$

- 相插值多项式：

$$\beta_j = \frac{1}{N} \sum_{k=0}^{N-1} f_k w^{-kj}, \quad j = 0, 1, \dots, N - 1$$

式中， $w = e^{i\frac{2\pi}{N}}$ 。

## 三 函数逼近

### 总结

Note: 这里为了让总结看起来更顺畅, 没有遵循原 PPT 和书中的符号规范, 转而应用了比较统一的符号。这些符号规范仅在总结一部分使用。

- 最佳平方逼近求解:

- 给定一组规范正交函数基时:

- \* 写出规范正交函数基  $\{\phi_m\}_{m=0}^n$
    - \* 计算系数  $a_m = \frac{(f, \phi_m)}{(\phi_m, \phi_m)}$
    - \* 写出逼近多项式  $\phi^* = \sum_{m=0}^n a_m \phi_m$

- 给定一组非规范正交函数基时:

- \* 写出非规范正交函数基  $\{\phi_m\}_{m=0}^n$
    - \* 构建法方程组:

$$m_{ij} = (\phi_j, \phi_i), \quad b_i = (f, \phi_i)$$

- \* 求解线性方程组  $\mathbf{Ma} = \mathbf{b}$ , 得到系数  $a_m$
    - \* 写出逼近多项式  $\phi^* = \sum_{m=0}^n a_m \phi_m$

- Legendre 多项式作最佳平方逼近的收敛速度: 高阶时控制在  $1/\sqrt{n}$  以下

- 最小二乘逼近求解: 计算过程可视作最佳平方逼近的离散形式

- 给定基底  $\{\phi_m\}_{m=0}^n$  和采样点  $\{x_k\}_{k=0}^N$
  - 估计一个误差系数  $\rho(x_i)$
  - 用半内积构建法方程组:

$$m_{ij} = \sum_{k=0}^N \rho(x_k) \phi_j(x_k) \phi_i(x_k), \quad b_i = \sum_{k=0}^N \rho(x_k) f(x_k) \phi_i(x_k)$$

- 求解线性方程组  $\mathbf{M}\mathbf{a} = \mathbf{b}$ , 得到系数  $a_m$
- 写出逼近多项式  $\phi^* = \sum_{m=0}^n a_m \phi_m$
- 一致逼近求解:
  - 给定原函数  $f \in C[a, b]$  和基底  $\{\phi_m\}_{m=0}^n$
  - 写出待定逼近多项式  $\phi(x) = \sum_{m=0}^n a_m \phi_m(x)$
  - 写出误差函数  $E(x) = f(x) - \phi(x)$
  - 根据切比雪夫交错点组定理, 写出求解条件:
    - \* 交错条件:  $E(x_i) = -e(x_{i-1})$
    - \* 偏差点条件: 除端点作为偏差点的情况, 必有  $f'(x_i) = p'_n(x_i)$
  - 由偏差点条件可以解出一系列偏差点  $x_i(a)$ , 再代入交错条件中, 解出系数  $a_m$

## (一) 通用理论

### 1. 问题模型

**待填写: (问题) 函数逼近问题的数学模型**

最佳逼近问题: 找  $\phi^* \in \Phi$ , 使得

$$\|f - \phi^*\| = \min_{\phi \in \Phi} \|f - \phi\|$$

### 2. 逼近准则

(1) 最小二乘准则 (L 范数)

(2) 一致逼近准则 ( $L^\infty$  范数)

### 3. 核心定理

## (二) 具体逼近方法

### 1. 最优平方逼近

#### (1) 基础理论

问题: 给定一个线性子空间  $\Phi$ , 找  $\phi^* \in \Phi$ , 使得

$$\|f - \phi^*\|_2 = \min_{\phi \in \Phi} \|f - \phi\|_2$$

问题等价于

$$\|f - \phi^*\|_2^2 = \min_{\phi \in \Phi} \|f - \phi\|_2^2$$

于是构造出一个辅助函数：

$$\begin{aligned} I &= \|f - \phi\|_2^2 = \left( f - \sum_{i=0}^n a_i \phi_i, f - \sum_{i=0}^n a_i \phi_i \right) \\ &= \sum_{i=0}^n \sum_{j=0}^n a_i a_j (\phi_i, \phi_j) - 2 \sum_{i=0}^n a_i (f, \phi_i) + (f, f) \end{aligned}$$

式中， $\{\phi_i\}$  为  $\Phi$  的一组基。

### a. 法方程

多元函数取得最小值的条件：对各变量偏导为 0、在这一点的二阶偏导数构成的矩阵正定。

各变量偏导为 0 即导出法方程：

$$\sum_{j=0}^n a_j (\phi_j, \phi_i) = (f, \phi_i), \quad i = 0, 1, \dots, n$$

或写成矩阵形式，

$$\begin{bmatrix} (\phi_0, \phi_0) & (\phi_1, \phi_0) & \cdots & (\phi_n, \phi_0) \\ (\phi_0, \phi_1) & (\phi_1, \phi_1) & \cdots & (\phi_n, \phi_1) \\ \vdots & \vdots & \ddots & \vdots \\ (\phi_0, \phi_n) & (\phi_1, \phi_n) & \cdots & (\phi_n, \phi_n) \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} (f, \phi_0) \\ (f, \phi_1) \\ \vdots \\ (f, \phi_n) \end{bmatrix}$$

**定理 (最小二乘逼近法方程).** 设  $\Phi$  为线性子空间， $\{\phi_0, \phi_1, \dots, \phi_n\}$  为  $\Phi$  的一组基，则函数  $f$  在  $\Phi$  中的最优平方逼近  $\phi^*$  可表示为

$$\phi^* = \sum_{j=0}^n a_j^* \phi_j$$

其中，系数  $a_j^*$  满足法方程组

$$\sum_{j=0}^n a_j (\phi_j, \phi_i) = (f, \phi_i), \quad i = 0, 1, \dots, n$$

### b. 最佳平方逼近的性质

**引理 (最佳平方逼近的正交性质).** 设  $\phi^*$  为函数  $f$  在子空间  $\Phi$  中的最佳平方逼近，则对任意  $\phi \in \Phi$ ，都有

$$f - \phi^* \perp \Phi$$

**推论 (最佳平方逼近的勾股定理).** 设  $\phi^*$  为函数  $f$  在子空间  $\Phi$  中的最佳平方逼近，则有

$$\|f\|_2^2 = \|f - \phi^*\|_2^2 + \|\phi^*\|_2^2$$

## (2) 幂函数基的逼近

考虑连续函数在  $n$ -次多项式空间  $P_n$  中的最佳平方逼近问题。

若选取  $\{x^m\}_{m=0}^n$  作为基函数，在  $[0, 1]$  区间上求解最佳平方逼近，则法方程组的系数矩阵为 Hilbert 矩阵：

**定义 (Hilbert 矩阵).** 阶数为  $n + 1$  的 Hilbert 矩阵定义为

$$H_{ij} = \frac{1}{i+j-1}, \quad i, j = 1, 2, \dots, n+1$$

即

$$H_n = \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} & \cdots & \frac{1}{n+1} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \cdots & \frac{1}{n+2} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} & \cdots & \frac{1}{n+3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{1}{n+1} & \frac{1}{n+2} & \frac{1}{n+3} & \cdots & \frac{1}{2n+1} \end{bmatrix}$$

这一逼近形式的问题：

- Hilbert 矩阵病态，随着  $n$  的增大，条件数迅速增大，导致数值解不稳定
- 幂函数基不正交，导致法方程系数矩阵接近奇异

## (3) 正交多项式逼近

### a. 广义傅里叶展开

在  $C[a, b]$  中，取一个规范正交函数组  $\{\phi_m\}_{m=0}^n$ ，即确定了一个用于逼近的线性子空间  $\Phi_n = \text{span}\{\phi_0, \phi_1, \dots, \phi_n\}$ 。

相应地，对于任意给定函数  $f$ ，存在唯一的最佳平方逼近  $\phi^* \in \Phi_n$ ，且

$$\phi^* = \sum_{i=0}^n a_i^* \phi_i, \quad a_i^* = (f, \phi_i)$$

**定义 (广义傅里叶展开).** 设  $\{\phi_m\}_{m=0}^\infty$  为  $C[a, b]$  中的规范正交函数组，则对任意  $f \in C[a, b]$ ，都有

$$f \sim f_\infty = \sum_{i=0}^\infty a_i \phi_i, \quad a_i = (f, \phi_i)$$

$f_\infty$  称为函数  $f$  在  $\{\phi_m\}$  下的广义傅里叶展开。

**定理 (广义傅里叶展开的收敛性).** 设  $\{\phi_m\}_{m=0}^\infty$  为  $C[a, b]$  中的规范正交函数组，则对任意  $f \in C[a, b]$ ，其广义傅里叶展开在  $L_2$  范数下收敛于  $f$ ，即

$$\lim_{n \rightarrow \infty} \|f - f_n\|_2 = 0$$

其中， $f_n = \sum_{i=0}^n a_i \phi_i$ 。

### b. Legendre 多项式作最佳平方逼近

设内积为  $[-1, 1]$  上的权 1 内积, 给定  $f \in C[-1, 1]$ , 则  $f$  在 Legendre 多项式空间  $P_n$  中的最佳平方逼近为

$$I_n f(x) = \sum_{k=0}^n \frac{(f, \phi_k)}{\|\phi_k\|^2} \phi_k(x)$$

性质:

- 收敛性:  $I_n f \rightarrow f, n \rightarrow \infty$ , 且收敛到  $L^2((-1, 1))$  空间中
- 收敛速度估计: 若  $f \in C^2[-1, 1]$ , 则  $\forall \epsilon > 0, \exists N >$ , 使得

$$\|f - I_n f\|_\infty \leq \frac{\epsilon}{\sqrt{n}}, \quad n \geq N$$

即对于  $n$  次的 Legendre 多项式逼近, 误差会被控制在  $O(\frac{1}{\sqrt{n}})$  范围内

### (4) Legendre 多项式的零平方误差最小性质

**定理** (Legendre 多项式的零平方误差最小性质). 设  $f \in C[-1, 1]$ , 则在所有满足  $\deg(p_n) \leq n$  且  $(p_n, x^k) = 0, k = 0, 1, \dots, n-1$  的多项式中, Legendre 多项式  $P_n(x)$  使得平方误差  $\|f - P_n\|_2$  最小。

## 2. 最小二乘逼近: 最优平方逼近的离散化形式

待填写: (问题) 最小二乘逼近的问题模型: 确定的输入-输出关系叠加系统噪声和随机误差, 希望拟合输入-输出关系的具体形式

### (1) 基础理论

核心假设: 观测数据  $f_i$  和正确输  $\phi^*(x_i)$  之间的误差  $a_i \xi_i$  是一个独立同分布的零均值随机变量。

假设导出: 当拟合结果  $\phi$  取得正确的关系  $\phi = \phi^*$  时, (归一化) 误差的方差最小。

数学化:  $\phi^*$  是以下优化问题的极小解:

$$\phi^* = \arg \min_{\phi \in \Phi} \sigma^2 = \arg \min_{\phi \in \Phi} \frac{1}{m+1} \sum_{i=0}^m [f_i - \phi(x_i)]^2$$

### (2) 数学方法: “半” 内积

**定理** (Haar 条件). 若给定一组采样点  $\{x_i\}_{i=0}^m$ , 且  $\Phi$  中的任意非 0 元素在这组采样点上的采样结果不全为 0, 则称这组采样点满足 Haar 条件。

满足 Haar 条件时, 可以利用在采样点上的采样结果, 定义一个“半”内积:

$$(f, g)_m = \sum_{i=0}^m \frac{f(x_i)g(x_i)}{a_i^2}$$

此时这个“半”内积在  $\Phi$  上是一个真正的内积。因此上述问题等价于求解法方程  $AU=b$ :

$$\begin{bmatrix} (\phi_0, \phi_0)_m & (\phi_1, \phi_0)_m & \cdots & (\phi_n, \phi_0)_m \\ (\phi_0, \phi_1)_m & (\phi_1, \phi_1)_m & \cdots & (\phi_n, \phi_1)_m \\ \vdots & \vdots & \ddots & \vdots \\ (\phi_0, \phi_n)_m & (\phi_1, \phi_n)_m & \cdots & (\phi_n, \phi_n)_m \end{bmatrix} \begin{bmatrix} u_0 \\ u_1 \\ \vdots \\ u_n \end{bmatrix} = \begin{bmatrix} (f, \phi_0)_m \\ (f, \phi_1)_m \\ \vdots \\ (f, \phi_n)_m \end{bmatrix}$$

Note: 我们注意到，在求内积过程中，分母上表示采样点随机误差幅度的  $a_i^2$  仍然是不确定的，此时需要基于经验或先验知识给出估计。常用的估计有  $a_i \propto 1$ 、 $a_i \propto f_i$ 、 $a_i \propto x_i$ 、 $a_i \propto x_i^2$  等。

### 3. 最佳一致逼近

**待填写：(问题) 最佳一致逼近的数学描述**

最佳一致逼近中用到一些符号：

- $\Delta(f, p_n) = \|f - p_n\|_\infty$  称为  $p_n$  关于  $f$  的偏差
- $E_n = \inf_{p_n \in P_n} \Delta(f, p_n)$  称为  $f$  在  $P_n$  中的最小偏差
- 若  $f(x_i) - p_n(x_i) = \sigma \Delta(f, p_n)$ ,  $\sigma = \pm 1$ , 则称  $x_i$  为  $p_n$  关于  $f$  的偏差点。
  - 若  $\sigma = 1$ , 则称  $x_i$  为正偏差点
  - 若  $\sigma = -1$ , 则称  $x_i$  为负偏差点

#### (1) 偏差泛函与最优逼近多项式存在性定理

假设给定一个逼近多项式  $p_n = a_0 + a_1x + a_2x^2 + \dots + a_nx^n$ , 则可定义偏差泛函:

$$\phi(f, \mathbf{a}) = \|f - p_n\|_\infty = \|f - (a_0 + a_1x + \dots + a_nx^n)\|_\infty$$

这个泛函具有以下性质：

- 连续性:  $\phi$  对  $\mathbf{a}$  连续
- 对  $\mathbf{a}$  下凸:  $\forall \mathbf{a}_1, \mathbf{a}_2 \in \mathbb{R}^{n+1}$ ,  $\forall t \in [0, 1]$ , 都有

$$\phi(f, t\mathbf{a}_1 + (1-t)\mathbf{a}_2) \leq t\phi(f, \mathbf{a}_1) + (1-t)\phi(f, \mathbf{a}_2)$$

- 正定性: 对于任意  $\mathbf{a} \neq \mathbf{0}$ , 都有  $\phi(f, \mathbf{a}) > 0$

正定性等价于:  $\phi(0; \mathbf{a})$  在单位球面  $\sum_{i=0}^n a_i^2 = 1$  上有正的最小值  $\mu$

利用偏差泛函的性质, 可以证明最佳一致逼近多项式的存在性:

**定理 (最佳一致逼近存在性定理).** 设  $f \in C[a, b]$ , 则在  $n$  次多项式空间  $P_n$  中, 存在一个多项式  $p_n^* \in P_n$ , 使得

$$\|f - p_n^*\|_\infty = \min_{p_n \in P_n} \|f - p_n\|_\infty$$

即  $P_n$  中关于  $f \in C[a, b]$  的最小偏差是可以达到的。

## (2) 最佳一致逼近多项式的性质

**定义** (Chebyshev 交错点组). 若  $x_1 < x_2 < \dots < x_m$  是  $p_n$  关于  $f$  的轮流为正负的偏差点, 则称其为一个 Chebyshev 交错点组。

**定理** (最佳一致逼近多项式的 Chebyshev 交错定理).  $p_n^*$  是  $f$  的最佳一致逼近多项式的充要条件是存在一个元素个数为  $n+2$  的 Chebyshev 交错点组  $x_0 < x_1 < \dots < x_{n+1}$

**推论.** 最佳一致逼近多项式是一个 Lagrange 插值多项式, 其插值节点在  $(a, b)$  内。

**定理** (最佳一致逼近多项式的唯一性定理). 设  $f \in C[a, b]$ , 则  $f$  在  $n$  次多项式空间  $P_n$  中的最佳一致逼近多项式  $p_n^*$  是唯一的。

## (3) Chebyshev 多项式的零一致误差最小性质

**定理** (Chebyshev 多项式的零一致误差最小性质). 设  $T_n(x)$  为  $n$  次 Chebyshev 多项式, 则在  $[-1, 1]$  上, 任意单位首项系数的  $n$  次多项式  $p_n(x)$  均满足

$$\|p_n\|_\infty \geq \|T_n\|_\infty = \frac{1}{2^{n-1}}$$

即 Chebyshev 多项式在  $[-1, 1]$  上具有最小的无穷范数。

## (4) 交错点组求解方法

### a. 解析求解: 利用交错点组定理

$$\begin{aligned} E(x_i) &= -E(x_{i-1}), \quad \forall x_i \\ f'(x_i) &= p'_n(x_i), \quad \forall x_i \neq a, b \end{aligned}$$

### b. 数值求解: Remez 算法

- 选择初始交错点组  $\{x_i^{(0)}\}_{i=0}^{n+1}$
- 计算当前逼近多项式  $p_n^{(k)}$ , 并计算偏差  $E^{(k)}(x) = f(x) - p_n^{(k)}(x)$
- 找到新的交错点组  $\{x_i^{(k+1)}\}_{i=0}^{n+1}$ , 即  $E^{(k)}(x)$  的极值点
- 若  $\|E^{(k+1)}\|_\infty$  与  $\|E^{(k)}\|_\infty$  足够接近, 则停止迭代, 否则返回步骤 2



## 四 数值微积分

### 总结

机械积分公式:

$$I_n(f) = \sum_{k=0}^n A_k f(x_k)$$

### 数值积分

- 数值积分方法的代数精度

- 定义: 可以准确积分 n 次多项式, 不能准确积分 n+1 次
  - 判据: 当有 n 个积分节点时,
    - \* 以下公式确保至少达到 n-1 阶:

$$A_k = I(L_k) = \int_a^b L_k(x) \rho(x) dx$$

\* 最多达到 2n-1 阶: 当且仅当积分节点为权函数  $\rho(x)$  对应的正交多项式的 n 个不同实根时

- 数值积分方法的代数稳定性

- 定义: 对任意 n 阶的求积公式, 求积系数的绝对值之和有确定的上界
  - 判据: 对任意的阶数 n, 所有的求积系数均大于 0, 则该方法是一致稳定的
  - Newton-Cotes 公式不是一致稳定的, 但复合积分方法和 Gauss 求积公式是一致稳定的。
- Newton-Cotes 公式: 用多项式函数插值近似函数, 再对插值多项式积分
- Newton-Cotes 公式的具体应用
  - 闭型: Newton-Cotes 公式系数和余项表
  - 开型: 中点公式和余项
- 复合求积方法 (以下规定  $h_k = x_{k+1} - x_k$ , 且  $h_{-1} = h_n = 0$ )

- 复合中点公式:  $A_k = h_k$
- 复合梯形公式:  $A_k = \frac{h_{k-1} + h_k}{2}$
- 复合 Simpson 公式: 即  $S_{2n} = \frac{1}{3}T_n + \frac{2}{3}H_n = \frac{4T_{2n} - T_n}{3}$
- 外推与 Romberg 方法

– Neville 递推式: 需记忆!

$$T_{ik} = T_{i,k-1} + \frac{T_{i,k-1} - T_{i-1,k-1}}{\frac{h_{i-k}^2}{h_i^2} - 1}$$

– Romberg 方法: Neville 在二分加密过程的特殊情况需记忆!

$$T_{ik} = T_{i,k-1} + \frac{T_{i,k-1} - T_{i-1,k-1}}{4^k - 1}$$

- Gauss 求积公式: 在正交多项式零点处取积分节点, 且系数为  $A_k^{(n)} = I(L_k^{(n)})$ , 即第 k 个节点处的采样系数为该节点处的 Lagrange 基函数的带权积分

## 数值微分

- 基础方法: 差商近似导数
- 插值型数值微分:

– 线性插值:

$$\begin{aligned} p_1(x) &= f[x_0, x_1]x + f[x_0] \\ p'_1(x) &= f[x_0, x_1] \end{aligned}$$

– 抛物插值:

$$\begin{aligned} p_2(x) &= f[x_0] + f[x_0, x_1](x - x_0) + f[x_0, x_1, x_2](x - x_0)(x - x_1) \\ p'_2(x) &= f[x_0, x_1] + f[x_0, x_1, x_2](2x - x_0 - x_1) \\ p''_2(x) &= 2f[x_0, x_1, x_2] \end{aligned}$$

- 外推方法: 等距节点、二分加密下, 按照

$$\begin{aligned} G_1(h) &= G(h) \\ G_{m+1}(h) &= G_m(h/2) + \frac{G_m(h/2) - G_m(h)}{2^m - 1} = \frac{4^m G_m(h/2) - G_m(h)}{4^m - 1} \end{aligned}$$

- 误差分析:

- 基础方法：向前、向后差商误差为  $O(h)$ ，中心差商误差为  $O(h^2)$ 。此外，中心差商求二阶导也可以达到  $O(h^2)$
- 插值型数值微分：n 次插值多项式的导数误差为  $O(h^{n+1-k})$ ，其中 k 为求导阶数。此外，抛物插值在中心点处的二阶导数精度可以加一阶。
- 外推方法：每次外推使误差阶数增加 2

## (一) 数值积分

### 1. 通用理论

#### (1) 积分问题模型

**待填写：(问题) 关注的对象是带权积分**

- 目标：构造一种不依赖于函数  $f$  具体表达形式的近似积分方法
- 基本思想：用简单函数近似被积函数，再计算简单函数的准确积分

#### (2) 求积公式核心

#### (3) 代数精度

定义. 如果

$$\tilde{I}(f) = I(f), \quad \forall f \in P_n$$

则称  $\tilde{I}$  的代数精度至少为  $n$ 。进一步，若存在一个  $P_{n+1}$  中的多项式  $f$ ，使得  $\tilde{I}(f) \neq I(f)$ ，则称  $\tilde{I}$  的代数精度为  $n$ 。

定理 (至少  $n-1$  阶代数精度的判据). 设  $L_k(x)$  为以  $x_1, \dots, x_n$  为节点时的第  $k$  个插值基函数，则当且仅当

$$A_k = I(L_k) = \int_a^b L_k(x) \rho(x) dx$$

时， $I_n(\cdot)$  具有至少  $n-1$  阶代数精度。

定理 (n 个求积节点的最大可能代数精度). n 个求积节点上的求积公式最大可能代数精度为  $2n-1$ ，当且仅当求积节点恰为权函数  $\rho(x)$  对应的正交多项式的 n 个不同实根时，达到这一代数精度。

#### (4) 一致稳定性

设  $\tilde{I}_n$  是一个线性积分法，n 为渐进参数（对应自由度个数）。

定义 (一致稳定性). 设对于一个方法导出的  $n$  阶机械求积公式为

$$\tilde{I}_n(f) = \sum_{k=0}^n A_k^{(n)} f(x_k)$$

若存在一个  $M > 0$ , 使得对任意  $n$  阶的求积公式, 均能满足

$$\sum_{k=0}^n |A_k^{(n)}| \leq M$$

则称该数值积分法是一致稳定的。

**定理.** 若某个积分法对任意的阶数  $n$  和序数  $k$ , 均有  $A_k^{(n)} > 0$ , 则这一积分法是一致稳定的。

考虑一致稳定性的意义在于, 当数据存在误差时, 有偏数据积分结果的误差可以由一直稳定性的参数  $M$  和误差的上界  $\delta$  的乘积控制。

- Simpson 公式是一致稳定的。
- Gauss 求积公式总是一致稳定的。

## 2. 具体求积方法

### (1) Newton-Cotes 公式

#### a. 导出思路

用多项式函数整体插值近似求积区间内的函数值, 再对插值多项式进行积分, 得到近似积分公式。即如下过程导出:

- 选取插值节点  $\{x_i\}_{i=0}^n$ 。注意, 这些节点不一定全在积分区间内。
- 构造  $n$ -次 Lagrange 插值多项式

$$P_n(x) = \sum_{i=0}^n f(x_i) l_i(x)$$

其中,  $l_i(x)$  为 Lagrange 基函数:

$$l_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x - x_j}{x_i - x_j}$$

- 对插值多项式进行积分, 得到近似积分公式

$$\tilde{I}(f) = I(P_n) = \sum_{k=0}^n f(x_k) I(L_k) = \sum_{k=0}^n A_k^{(n)} f(x_k)$$

**b. 性质:**

- 近似积分的精度与插值多项式逼近精度相关
- 积分法只与函数  $f$  在插值节点处的取值有关

**c. (闭型) Newton-Cotes 公式**

取  $x_i$  为等距节点, 且  $a = x_0 < x_1 < \dots < x_n = b$ , 则称所得求积公式为闭型 Newton-Cotes 公式。积分系数为:

$$A_k^{(n)} = (b - a)C_k^{(n)}$$

(需记忆) Newton-Cotes 公式系数表:

n	weighted factor					
1	1	1				
2	1	4	1			
3	1	3	3	1		
4	7	32	12	32	7	

数值积分结果

$$I_n(f) = (b - a) \sum_{k=0}^n C_k^{(n)} f(x_k)$$

低阶公式:

- 一阶: 梯形公式:

$$\tilde{I}_1(f) = \frac{b-a}{2} [f(a) + f(b)]$$

$$R_1(x) = f(x) - P_1(x) = f[a, b, x](x-a)(x-b)$$

$$I(f) - I_1(f) = \int_a^b R_1(x) dx = \frac{f''(\xi)}{2!} \int_a^b (x-a)(x-b) dx = -\frac{(b-a)^3}{12} f''(\xi)$$

- 二阶: Simpson 公式:

$$\tilde{I}_2(f) = \frac{b-a}{6} [f(x_0) + 4f(x_1) + f(x_2)]$$

$$I(f) - \tilde{I}_2(f) = -\frac{1}{90} \left( \frac{b-a}{2} \right)^5 f^{(4)}(\xi)$$

注意: Simpson 公式的误差是 4 阶导数, 而不是 3 阶。

- 三阶: 3/8 - 规则:

$$\tilde{I}_3(f) = \frac{b-a}{8} [f(x_0) + 3f(x_1) + 3f(x_2) + f(x_3)]$$

$$I(f) - \tilde{I}_3(f) = -\frac{3}{80} \left(\frac{b-a}{3}\right)^5 f^{(4)}(\xi)$$

- 奇数阶的 Newton-Cotes 公式的代数精确度是 n; 偶数阶的 Newton-Cotes 公式的代数精确度是 n+1。这是由于取等距节点插值时, 将均差转化为余项后, 误差中的多项式在中点之前的误差和中点之后的误差恰好抵消。

#### d. 开型 Newton-Cotes 公式

取等距节点  $a = x_{-1} < x_0 < x_1 < \dots < x_n < x_{n+1} = b$ , 仅使用中间的  $x_0 < x_1 < \dots < x_n$  作为插值节点, 则称所得求积公式为开型 Newton-Cotes 公式。

$$I(f) = (b-a) \sum_{k=0}^n b_k(n) f(x_k) + \bar{R}_n(f)$$

$$b_k^{(n)} = \frac{1}{b-a} \int_a^b L_k(x) dx = \frac{1}{b-a} \int_a^b \prod_{\substack{j=0 \\ j \neq k}}^n \frac{x-x_j}{x_k-x_j} dx$$

$n=0$  的开型 Newton-Cotes 公式称为中点公式:

$$I(f) = (b-a)f(x_0) + \bar{R}_0(f), \quad x_0 = \frac{a+b}{2}$$

中点公式的误差:

$$\bar{R}_0[f] = \int_a^b f[x_0, x](x-x_0) dx = \frac{1}{3} h^3 f''(\xi), \quad h = \frac{b-a}{2}$$

### (2) Hermite 积分法

类似 Newton-Cotes 公式的思路, 采用更高阶的近似:

在区间  $[a, b]$  上, 用 f 的三次 Hermite 插值近似  $f(x)$ , 再对插值多项式进行积分, 即可得到 Hermite 积分公式。

仅在区间端点上采样时, 得到下式:

$$I(f) = \frac{b-a}{2} [f(a) + f(b)] + \frac{(b-a)^2}{12} [f'(a) - f'(b)] + \frac{1}{720} (b-a)^5 f^{(4)}(\xi)$$

### (3) 复合求积方法

在区间内取 n 个采样点, 不要求均匀, 即取

$$a = x_0 < x_1 < \dots < x_n = b$$

并记第  $k$  个区间长度为下一个点减去这一个点:

$$h_k = x_{k+1} - x_k, \quad k = 0, 1, \dots, n-1$$

同时, 为了统一表示形式, 约定  $h_{-1} = h_n = 0$ 。

### a. 复合中点公式

$$I(f) = \sum_{k=0}^{n-1} h_k f\left(\frac{x_k + x_{k+1}}{2}\right) + R_n(f)$$

机械求积公式形式:

$$H_n(f) = \sum_{k=0}^n A_k^{(n)} f(x_k)$$

$$A_k = h_k$$

### b. 复合梯形公式

$$I(f) = \sum_{k=0}^{n-1} \frac{h_k}{2} (f(x_k) + f(x_{k+1})) + R_n(f)$$

机械求积公式形式 (约定  $h_{-1} = h_n = 0$ ):

$$\begin{aligned} T_n(f) &= \sum_{k=0}^n A_k^{(n)} f(x_k) \\ A_k &= \frac{h_{k-1} + h_k}{2} \end{aligned}$$

误差分析:

$$\begin{aligned} R_n(f) &= -\frac{1}{12} \sum_{k=0}^{n-1} h_k^3 f''(\xi_k), \quad \xi_k \in [x_k, x_{k+1}] \\ &= -\frac{1}{12} \sum_{k=0}^{n-1} h_k^3 f''(\eta), \quad \eta \in [a, b] \end{aligned}$$

设  $h = \max_{0 \leq k \leq n-1} h_k$ , 则

$$|R_n(f)| \leq \frac{b-a}{12} h^2 |f''(\eta)|$$

递推关系:

$$T_{2n}(f) = \frac{T_n(f) + H_n(f)}{2}$$

即在每个积分区间上增加一个节点后, 复合梯形公式即为增加节点前的复合梯形公式与复合中点公式的平均值。

### c. 改造复合梯形公式

使用 Hermite 积分法，考虑等距节点的复合梯形公式，可以将复合梯形公式的积分精度改造到四阶：

$$I(f) = \sum_{k=0}^{n-1} \frac{h_k}{2} [f(x_k) + f(x_{k+1})] + \frac{h_k^2}{12} [f'(x_k) - f'(x_{k+1})] + \frac{1}{720} (b-a)^5 f^{(4)}(\xi)$$

机械求积公式：

$$\begin{aligned} \tilde{T}_n(f) &= \sum_0^n A_k^{(n)} f(x_k) \\ A_k &= \begin{cases} \frac{5}{12}h, & k = 0, n \\ \frac{13}{12}h, & k = 1, n-1 \\ h, & \text{其他} \end{cases} \end{aligned}$$

### d. 复合 Simpson 公式：一致稳定的求积公式

$$I(f) = \sum_{k=0}^{n-1} \frac{h_k}{6} [f(x_k) + 4f\left(\frac{x_k + x_{k+1}}{2}\right) + f(x_{k+1})] + R_n(f)$$

机械求积公式为

$$S_{2n}(f) = \frac{1}{3}T_n(f) + \frac{2}{3}H_n(f) = \frac{4T_{2n}(f) - T_n(f)}{3}$$

误差分析：

$$|R_{2n}(f)| \leq \frac{b-a}{2880} h^4 |f^{(4)}(\eta)|, \quad \eta \in [a, b]$$

其中  $h = \max_{0 \leq k \leq n-1} h_k$ 。

## (4) Romberg 加速方法

目的：在提升积分区间取点数时，充分利用利用已有的积分结果，通过外插法提升积分精度。

### a. 外插法

在取等距节点时，成立欧拉-麦克劳林公式：

定理 (Euler-Maclaurin 公式). 设  $f \in C^{2m+2}[a, b]$ , 则有

$$I(f) = T_n(f) + \sum_{k=1}^m \frac{B_{2k}}{(2k)!} h^{2k} [f^{(2k-1)}(b) - f^{(2k-1)}(a)] + R_{m+1}$$

其中,  $T_n(f)$  为复合梯形公式,  $h = \frac{b-a}{n}$ ,  $B_{2k}$  为 Bernoulli 数, 且

$$r_{m+1} = -\frac{B_{2m+2}}{(2m+2)!} (b-a) f^{(2m+2)}(\eta) h^{2m+2}, \quad \eta \in [a, b]$$

由 Euler-Maclaurin 公式可知 (只要将上式的级数移到左侧), 复合梯形公式可以写成关于  $h$  的渐进级数形式, 且 0 次项即为积分结果:

$$\begin{aligned} T_n(f) &\sim \tau_0 + \tau_1 h^2 + \tau_2 h^4 + \dots, \\ \tau_0 &= I(f) \end{aligned}$$

在逐步提升采样节点密度时, 我们总是采用等距节点, 则提升密度过程可以导出一个关于  $h$  的单调递减序列:

$$h_i = \frac{b-a}{n_i}, \quad n_i \text{ 单调递增}$$

于是可以考虑一个关于  $h^2$  的  $m$  次插值多项式 (具体  $m$  取决于精度要求):

$$\begin{aligned} \tilde{T}_{mm} &= a_0 + a_1 h^2 + a_2 h^4 + \dots + a_m h^{2m} \\ \tilde{T}_{mm}(h_i) &= T(h_i), \quad i = 0, 1, \dots, m \end{aligned}$$

于是将采样点的增加转化为了插值点的增加。

### b. Neville 算法: 外插法的计算

记  $\tilde{T}_{ik}$  为关于  $h^2$  的  $k$  次插值多项式, 满足的插值条件为最后的  $k+1$  个  $h_i$  采样点:

$$\tilde{T}_{ik}(h_j) = T(h_j), \quad j = i-k, i-k+1, \dots, i.$$

则可以通过低一阶的插值多项式构造高一阶的插值多项式:

$$\tilde{T}_{ik}(h) = \frac{(h^2 - h_{i-k})^2 \tilde{T}_{i,k-1}(h) + (h_i^2 - h^2) \tilde{T}_{i-1,k-1}(h)}{h_i^2 - h_{i-k}^2}$$

我们仅关注这一递推关系在  $h = 0$  处的取值, 于是得到 (重要递推关系, 需记忆!)

$$T_{ik} = T_{i,k-1} + \frac{\frac{T_{i,k-1} - T_{i-1,k-1}}{h_{i-k}^2} - 1}{\frac{h_i^2}{h_{i-k}^2} - 1}$$

插值多项式在 0 处的值  $T_{ik}$  有一个需要注意的性质:

- $T_{i0} = T_{n_i}$ , 即 0 阶插值多项式即为对应  $h$  的复合梯形公式积分结果, 或取  $n_i$  个等距节点时的复合梯形公式积分结果

### c. Romberg 方法

对于最简单的二分加密过程，有  $n_i = 2^i$ ，则  $h_i = \frac{b-a}{2^i}$ 。

于是递推式化为

$$T_{ik} = T_{i,k-1} + \frac{T_{i,k-1} - T_{i-1,k-1}}{4^k - 1}$$

特别地，当  $k=1$  时，有：

$$T_{i1} = \frac{4T_{i0} - T_{i-1,0}}{3} = \frac{4T_{2n} - T_n}{3}$$

即为 Simpson 公式。

**Romberg 方法的列表求解见例题。逐步外推过程。**

#### D. Romberg 方法的代数精度

若用 Romberg 方法计算出一个  $T_{mm}$ ，并将其作为积分结果，则该积分方法的代数精度至少为  $2m+1$ 。

## (5) Gauss 型求积

**定义.** 若  $x_1, \dots, x_n$  取为权函数  $\rho(x)$  对应的正交多项式  $\phi_n(x)$  的  $n$  个不同实根，则所得求积公式

$$I_n(f) = \sum_{k=1}^n A_k^{(n)} f(x_k), \quad A_k^{(n)} = I(L_k^{(n)}) = \int_a^b L_k^{(n)}(x) \rho(x) dx$$

有  $2n-1$  阶代数精度，相应的公式称为 *Gauss* 公式， $x_k$  称为 *Gauss* 点。

### a. Gauss 积分公式的基本性质

- Gauss 积分公式是一致稳定的
- Gauss 积分公式是收敛的，即当采样点个数趋于无穷大时，积分结果趋于真实值

### b. 设计 Gauss 积分公式的步骤

- 求正交多项式  $\phi_n$
- 求正交多项式的  $n$  个不同实根，作为 Gauss 点
- 计算积分系数  $A_k = I(L_k)$
- 得到求积公式

$$I_n(f) = \sum_{k=1}^n A_k f(x_k)$$

### c. Gauss 积分公式的应用

应用优势区：函数光滑性足够好。如果函数光滑性不太好，可根据奇点分片。

两种典型积分：

- Gauss-Legendre 求积：适用于  $\rho(x) = 1$  的情形

$$I(f) = \int_{-1}^1 f(x) dx$$

- Gauss-Chebyshev 求积：适用于  $\rho(x) = \frac{1}{\sqrt{1-x^2}}$  的情形

$$I(f) = \int_{-1}^1 \frac{f(x)}{\sqrt{1-x^2}} dx$$

在一般积分区间上，可以通过线性变换将积分区间映射到  $[-1, 1]$  上：

$$I(f) = \int_a^b f(x) dx = \frac{b-a}{2} \int_{-1}^1 f\left(\frac{b-a}{2}t + \frac{a+b}{2}\right) dt$$

## (6) 特殊积分处理

### a. 奇异积分

常用方法：

- 变量替换消奇点：如  $\int_0^1 x^{-1/n} f(x) dx = n \int_0^1 f(t^n) t^{n-2} dt$
- 奇性分离：将函数分为解析可积的有奇性部分和一个光滑部分，解析处理奇性部分，数值处理光滑部分
- 无穷区间：在 Laguerre 多项式零点或者 Hermite 多项式零点进行 Gauss 积分

### b. 振荡积分

对于形如

$$\int_a^b f(x) e^{i\omega x} dx$$

的问题，常用两种方法处理：

- 将问题区域分片，用高精度的分片多项式插值函数  $g$  逼近  $f$ ，再（解析或者半解析地）准确计算  $g$  的积分
- 若  $f$  为足够光滑的周期函数，则可用高次三角函数近似

## (7) 高维积分

### a. 蒙特卡洛方法

## (二) 数值微分

### 1. 基础方法

#### (1) 向前差商: 误差 $O(h)$

$$f'(x) = \frac{f(x+h) - f(x)}{h} + O(h) = f[x, x+h] + O(h)$$

#### (2) 向后差商: 误差 $O(h)$

$$f'(x) = \frac{f(x) - f(x-h)}{h} + O(h) = f[x-h, x] + O(h)$$

#### (3) 中心差商: 误差 $O(h^2)$

$$f'(x) = \frac{f(x+h) - f(x-h)}{2h} + O(h^2) = f[x-h, x+h] + O(h^2)$$

高阶导数:

$$f''(x) = \frac{f(x+h) - 2f(x) + f(x-h)}{h^2} + O(h^2) = 2f[x-h, x, x+h] + O(h^2)$$

### 2. 高精度方法

#### (1) 插值型数值微分

思路: 设函数  $g$  是满足在  $\{x_k\}_{k=0}^n$  上插值条件的  $f$  的近似函数, 插值型求导方法即令

$$f^{(n)}(x) \approx g^{(n)}(x)$$

作为  $f$  的导数的近似。

误差分析:

$$\begin{aligned} f(x) - p_n(x) &= f[x_0, x_1, \dots, x_n, x] w_{n+1}(x) \\ w_{n+1}(x) &= (x - x_0)(x - x_1) \dots (x - x_n) \end{aligned}$$

求导可得

$$\begin{aligned} f'(x_k) - p'_n(x_k) &= f[x_0, x_1, \dots, x_n, x_k] w'_{n+1}(x_k) \\ f''(x_k) - p''_n(x_k) &= 2f[x_0, x_1, \dots, x_n, x_k, x_k] w'_{n+1}(x_k) + f[x_0, x_1, \dots, x_n, x_k] w''_{n+1}(x_k) \end{aligned}$$

即通常情况下，一阶导有  $O(h^n)$  误差，二阶导有  $O(h^{n-1})$  误差。

### a. 线性插值公式

$$\begin{aligned} p_1(x) &= f[x_0] + f[x_0, x_1](x - x_0) \\ p'_1(x) &= f[x_0, x_1] \end{aligned}$$

### b. 抛物插值

$$\begin{aligned} p_2(x) &= f[x_0] + f[x_0, x_1](x - x_0) + f[x_0, x_1, x_2](x - x_0)(x - x_1) \\ p'_2(x) &= f[x_0, x_1] + f[x_0, x_1, x_2](2x - x_0 - x_1) \\ p''_2(x) &= 2f[x_0, x_1, x_2] \end{aligned}$$

等距节点时，有

$$\begin{aligned} f'(x_0) &= \frac{1}{2h}[-3f(x_0) + 4f(x_1) - f(x_2)] + O(h^2) \\ f'(x_1) &= \frac{1}{2h}[-f(x_0) + f(x_2)] + O(h^2) \\ f'(x_2) &= \frac{1}{2h}[f(x_0) - 4f(x_1) + 3f(x_2)] + O(h^2) \\ f''(x_0) &= \frac{1}{h^2}[f(x_0) - 2f(x_1) + f(x_2)] + O(h^2) \end{aligned}$$

Note: 由于抛物插值有一定对称性，中心点的二阶导精度可以加一阶。

## (2) Richardson 外推加速

类似数值积分的外推法，若有多个不同区间长度下的中点微分公式近似结果，则可进行数值外推进行二分加密时，公式如下：

$$\begin{aligned} G_1(h) &= G(h) \\ G_{m+1}(h) &= G_m(h/2) + \frac{G_m(h/2) - G_m(h)}{2^m - 1} = \frac{4^m G_m(h/2) - G_m(h)}{4^m - 1} \end{aligned}$$



# 五 非线性方程求根 (Nonlinear Equations)

求解  $f(x) = 0$ , 其中  $f : [a, b] \rightarrow \mathbb{R}$  为连续函数。

- 二分法: 基于介值定理, 线性收敛但可靠。
- 迭代法: 将问题转化为等价不动点问题  $x = \phi(x)$ 。利用迭代格式  $x^{(k+1)} = \phi(x^{(k)})$  得到。

## (一) 不动点迭代法 (Fixed Point Iteration)

### 1. 基本概念

- 不动点: 若  $\xi = \phi(\xi)$ , 则称  $\xi$  为  $\phi$  的不动点。
- 等价性:  $f(x) = 0 \iff x = \phi(x)$ 。例如取  $\phi(x) = x - cf(x)$ 。

### 2. 收敛理论

定理 (Brouwer 不动点定理). 设  $\phi : D \rightarrow \mathbb{R}^n$  是凸紧集  $D$  上的连续映射, 且  $\phi(D) \subseteq D$ 。则  $\phi$  在  $D$  内至少存在一个不动点。(注: 在一维情形下, 即若  $\phi : [a, b] \rightarrow [a, b]$  连续, 则必有不动点。)

定理 (压缩映射原理 / Banach Fixed Point Theorem). 设  $\phi : D \rightarrow D$  是完备度量空间上的压缩映射, 即存在常数  $0 \leq L < 1$ , 使得:

$$\|\phi(x) - \phi(y)\| \leq L\|x - y\|, \quad \forall x, y \in D$$

则:

1.  $\phi$  在  $D$  内存在唯一不动点  $x^*$ 。
2. 对任意初值  $x_0 \in D$ , 迭代序列  $x_{k+1} = \phi(x_k)$  均收敛于  $x^*$ 。
3. 误差估计:

$$\|x_k - x^*\| \leq \frac{L^k}{1-L} \|x_1 - x_0\|$$

说明. 压缩映射证明思路

1. 证明  $\{x_k\}$  是柯西序列:  $\|x_{k+m} - x_k\| \leq \frac{L^k}{1-L} \|x_1 - x_0\|$ 。

2. 利用空间完备性，得极限  $x^*$  存在。

3. 利用连续性证明  $x^*$  是不动点。

4. 利用反证法证明唯一性。

### 3. 局部收敛性定理

**定理.** 设  $x^*$  为  $\phi(x)$  的不动点，且  $\phi'(x)$  在  $x^*$  邻域内连续。若

$$|\phi'(x^*)| < 1$$

则不动点迭代在  $x^*$  的某个邻域内局部收敛。

**说明.** 由微分中值定理：

$$x_{k+1} - x^* = \phi(x_k) - \phi(x^*) = \phi'(\xi_k)(x_k - x^*)$$

其中  $\xi_k$  介于  $x_k$  和  $x^*$  之间。取绝对值：

$$|x_{k+1} - x^*| = |\phi'(\xi_k)| |x_k - x^*|$$

若  $|\phi'(x^*)| < 1$ ，由于  $\phi'$  连续，存在邻域使得  $|\phi'(x)| \leq L < 1$ 。于是  $|x_{k+1} - x^*| \leq L|x_k - x^*|$ ，即收敛。

### 4. 收敛阶 (Order of Convergence)

**定义.** 若迭代序列  $\{x_k\}$  收敛于  $x^*$ ，且

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - x^*|}{|x_k - x^*|^p} = C \neq 0$$

则称该迭代过程是  $p$  阶收敛的。

**定理 (收敛阶判定).** 设  $\phi(x)$  在  $x^*$  附近足够光滑。

- 若  $0 < |\phi'(x^*)| < 1$ ，则为线性收敛 ( $p = 1$ )。
- 若  $\phi'(x^*) = 0$  且  $\phi''(x^*) \neq 0$ ，则为平方收敛 ( $p = 2$ )。
- 一般地，若  $\phi^{(k)}(x^*) = 0 (k = 1, \dots, p-1)$  且  $\phi^{(p)}(x^*) \neq 0$ ，则为  $p$  阶收敛。

**说明.** 推导过程 (利用泰勒展开) 将  $\phi(x_k)$  在  $x^*$  处展开：

$$x_{k+1} = \phi(x_k) = \phi(x^*) + \phi'(x^*)(x_k - x^*) + \frac{\phi''(x^*)}{2!}(x_k - x^*)^2 + \dots$$

注意到  $\phi(x^*) = x^*$ ，代入误差  $e_k = x_k - x^*$ ：

$$e_{k+1} = \phi'(x^*)e_k + \frac{\phi''(x^*)}{2!}e_k^2 + \dots + \frac{\phi^{(p)}(x^*)}{p!}e_k^p + O(e_k^{p+1})$$

- 若  $\phi'(x^*) \neq 0$ ，则  $e_{k+1} \approx \phi'(x^*)e_k$ ，即线性收敛。
- 若  $\phi'(x^*) = 0$  但  $\phi''(x^*) \neq 0$ ，则  $e_{k+1} \approx \frac{\phi''(x^*)}{2}e_k^2$ ，即平方收敛。

## 5. 收敛域 (Convergence Domain)

定义 (收敛域). 对于给定的不动点  $x^*$ , 使得迭代序列  $x_{k+1} = \phi(x_k)$  收敛于  $x^*$  的所有初始点  $x_0$  的集合, 称为该不动点的收敛域。

### (1) 求解收敛域

理论上求精确的收敛域很难, 但通常我们需要找到一个保证收敛的区间。

1. 解不等式: 首先解不等式  $|\phi'(x)| < 1$ 。这将给出一个或多个开区间  $I$ 。在这些区间内,  $\phi$  是压缩的。
2. 验证封闭性: 选取上述区间中包含不动点  $x^*$  的最大区间  $(a, b)$ 。检查该区间是否满足:

$$\forall x \in [a, b], \quad \phi(x) \in [a, b]$$

如果是, 则  $[a, b]$  是一个收敛区间。

## (二) 牛顿法 (Newton's Method)

### 1. 构造

- 几何视角: 在当前点  $(x_k, f(x_k))$  作切线, 取切线与  $x$  轴交点为下一次迭代值  $x_{k+1}$ 。
- 泰勒展开视角: 将  $f(x)$  在  $x_k$  处展开, 忽略高阶项:

$$f(x) \approx f(x_k) + f'(x_k)(x - x_k) = 0 \implies x = x_k - \frac{f(x_k)}{f'(x_k)}$$

迭代公式:

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$$

### 2. 收敛性分析

特殊的不动点迭代,  $\phi(x) = x - \frac{f(x)}{f'(x)}$ 。

定理 (局部二阶收敛性). 设  $x^*$  是  $f(x) = 0$  的单根 (即  $f(x^*) = 0, f'(x^*) \neq 0$ ), 且  $f(x)$  在  $x^*$  邻域内二阶连续可微。则存在  $x^*$  的邻域, 使得对任意  $x_0$  在此邻域内, 牛顿法生成的序列  $\{x_k\}$  收敛于  $x^*$ , 且为二阶收敛。

说明. 证明

1. 考察迭代函数  $\phi(x) = x - \frac{f(x)}{f'(x)}$  的导数:

$$\phi'(x) = 1 - \frac{(f'(x))^2 - f(x)f''(x)}{(f'(x))^2} = \frac{f(x)f''(x)}{(f'(x))^2}$$

2. 在单根  $x^*$  处,  $f(x^*) = 0$ , 故  $\phi'(x^*) = 0$ 。

3. 根据不动点收敛阶判定定理, 由于  $\phi'(x^*) = 0$ , 只要  $\phi''(x^*) \neq 0$ , 收敛阶至少为 2。

**误差渐进式:** 利用泰勒展开可得

$$e_{k+1} \approx \frac{f''(x^*)}{2f'(x^*)} e_k^2 \implies \lim_{k \rightarrow \infty} \frac{|e_{k+1}|}{|e_k|^2} = \left| \frac{f''(x^*)}{2f'(x^*)} \right|$$

### 3. 重根情形修正

若  $x^*$  是  $f(x) = 0$  的  $m$  重根 ( $m > 1$ ), 则  $f(x) = (x - x^*)^m g(x)$ , 且  $g(x^*) \neq 0$ 。此时,  $f(x)$  的前  $m - 1$  阶导数为 0. 标准牛顿法退化为线性收敛, 收敛因子为  $1 - 1/m$ 。

**修正牛顿法:** 为了恢复二阶收敛, 可使用:

$$x_{k+1} = x_k - m \frac{f(x_k)}{f'(x_k)}$$

$m$  为不动点的重数。

另:

令  $u := \frac{f}{f'}$ , 对  $u$  使用牛顿迭代法。

## (三) 加速方法

### 1. 割线法 (Secant Method)

**构造思想:** 为了避免牛顿法中计算导数  $f'(x_k)$  的代价, 用差商来近似导数:

$$f'(x_k) \approx \frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}}$$

将此代入牛顿公式, 得到迭代公式:

$$x_{k+1} = x_k - f(x_k) \frac{x_k - x_{k-1}}{f(x_k) - f(x_{k-1})}$$

此方法需要两个初值  $x_0, x_1$ 。

**收敛性分析:**

**定理 (割线法收敛阶).** 若  $x^*$  是单根, 且  $f$  二阶连续可微, 初值充分接近  $x^*$ , 则割线法是超线性收敛的, 收敛阶为黄金分割比:

$$p = \frac{1 + \sqrt{5}}{2} \approx 1.618$$

**说明.** 误差方程推导 设误差  $e_k = x_k - x^*$ 。割线法的误差关系满足:

$$e_{k+1} \approx C e_k e_{k-1}$$

设  $|e_{k+1}| \sim |e_k|^p$ , 则  $|e_k| \sim |e_{k-1}|^p \implies |e_{k-1}| \sim |e_k|^{1/p}$ 。代入误差方程:

$$|e_k|^p \sim |e_k| \cdot |e_k|^{1/p} = |e_k|^{1+1/p}$$

比较指数:  $p = 1 + \frac{1}{p} \implies p^2 - p - 1 = 0$ 。解得正根  $p = \frac{1+\sqrt{5}}{2} \approx 1.618$ 。

## 2. Aitken $\Delta^2$ 加速

构造：假设序列  $\{x_k\}$  线性收敛于  $x^*$ ，且渐进误差常数为  $C$ ：

$$x_k - x^* \approx C(x_{k-1} - x^*)$$

我们有三个连续点  $x_k, x_{k+1}, x_{k+2}$ 。近似认为  $C$  在这几步是不变的：

$$\frac{x_{k+2} - x^*}{x_{k+1} - x^*} \approx \frac{x_{k+1} - x^*}{x_k - x^*} \approx C$$

消去  $C$ ，解出  $x^*$  的估计值  $\hat{x}_k$ ：

$$(x_{k+2} - x^*)(x_k - x^*) \approx (x_{k+1} - x^*)^2$$

解得 Aitken 加速公式：

$$\hat{x}_k = x_k - \frac{(x_{k+1} - x_k)^2}{x_{k+2} - 2x_{k+1} + x_k} = x_k - \frac{(\Delta x_k)^2}{\Delta^2 x_k}$$

即使  $\{x_k\}$  只是线性收敛， $\{\hat{x}_k\}$  的收敛速度通常会快得多（高阶收敛）。

## 3. Steffensen 方法

Aitken 加速实质上是原不动点迭代步骤完成后，进一步提高精度的处理。

Steffensen 方法则将二阶差分直接结合到不动点迭代过程中，构成一种无需导数的二阶收敛算法。

构造过程：对于不动点迭代  $x = \phi(x)$ ，每一步迭代如下：

1. 给定  $x_k$ ，计算两步辅助点：

$$y_k = \phi(x_k), \quad z_k = \phi(y_k)$$

2. 将  $x_k, y_k, z_k$  视为 Aitken 加速中的三个点，直接计算加速值作为  $x_{k+1}$ ：

$$x_{k+1} = x_k - \frac{(y_k - x_k)^2}{z_k - 2y_k + x_k}$$

其中分母  $z_k - 2y_k + x_k = \phi(\phi(x_k)) - 2\phi(x_k) + x_k$ 。

Steffensen 方法对应的迭代函数

本质上是将单步不动点迭代函数  $\phi$  换成了  $\psi(x)$ ：

$$\psi(x) = x - \frac{(\phi(x) - x)^2}{\phi(\phi(x)) - 2\phi(x) + x} = \frac{x\phi(\phi(x)) - [\phi(x)]^2}{\phi(\phi(x)) - 2\phi(x) + x}$$

即  $x_{k+1} = \psi(x_k)$ 。

收敛性分析：

定理. Steffensen 方法实际上等价于对函数  $g(x) = x - \phi(x)$  使用牛顿法的一个近似形式（用差商代替导数）。若  $\phi(x)$  三阶连续可微，且  $\phi'(x^*) \neq 1$ ，则 Steffensen 方法是二阶收敛的。

$$\lim_{k \rightarrow \infty} \frac{|e_{k+1}|}{|e_k|^2} = C \neq 0$$

不需要计算导数，却能达到牛顿法的二阶收敛速度。

## (四) 非线性方程组

### 1. 基础理论与向量值微积分

#### (1) 非线性方程组

求解  $\mathbf{F}(\mathbf{x}) = \mathbf{0}$ , 其中  $\mathbf{F}: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  为向量值函数,  $\mathbf{x} = (x_1, \dots, x_n)^T$ 。

$$\mathbf{F}(\mathbf{x}) = \begin{pmatrix} f_1(x_1, \dots, x_n) \\ \vdots \\ f_n(x_1, \dots, x_n) \end{pmatrix}$$

#### (2) 向量值函数的导数 (Jacobian Matrix)

定义 (Jacobian 矩阵). 设  $\mathbf{F}$  在  $\mathbf{x}$  处可微, 其导数 (Jacobian 矩阵) 定义为:

$$\mathbf{J}(\mathbf{x}) = \mathbf{F}'(\mathbf{x}) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \cdots & \frac{\partial f_n}{\partial x_n} \end{pmatrix} \in \mathbb{R}^{n \times n}$$

#### (3) 多元 Taylor 展开

定理 (多元 Taylor 公式). 设  $\mathbf{F}$  二阶连续可微, 则在  $\mathbf{x}$  的邻域内:

$$\mathbf{F}(\mathbf{x} + \mathbf{h}) = \mathbf{F}(\mathbf{x}) + \mathbf{J}(\mathbf{x})\mathbf{h} + O(||\mathbf{h}||^2)$$

其中线性主部  $\mathbf{J}(\mathbf{x})\mathbf{h}$  是  $\mathbf{F}$  在  $\mathbf{x}$  处的微分。

#### (4) 多元压缩映射原理

定理 (Banach Fixed Point Theorem in  $\mathbb{R}^n$ ). 设  $D \subset \mathbb{R}^n$  是闭集,  $: D \rightarrow D$  是压缩映射, 即存在  $0 \leq L < 1$  使得:

$$||(\mathbf{x}) - (\mathbf{y})|| \leq L||\mathbf{x} - \mathbf{y}||, \quad \forall \mathbf{x}, \mathbf{y} \in D$$

则 在  $D$  内存在唯一不动点  $\mathbf{x}^*$ 。

推论 (局部收敛判据). 设  $\mathbf{x}^*$  是不动点。若  $\rho((\mathbf{x}^*)) < 1$  (谱半径小于 1), 则不动点迭代  $\mathbf{x}_{k+1} = (\mathbf{x}_k)$  在  $\mathbf{x}^*$  邻域内局部收敛。

### 2. 多元牛顿法 (Newton's Method for Systems)

构造思想: 对  $\mathbf{F}(\mathbf{x}) = \mathbf{0}$  在  $\mathbf{x}_k$  处进行线性化 (Taylor 展开保留一项):

$$\mathbf{F}(\mathbf{x}) \approx \mathbf{F}(\mathbf{x}_k) + \mathbf{J}(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k) = \mathbf{0}$$

解出  $\mathbf{x}$  即为下一次迭代值  $\mathbf{x}_{k+1}$ 。

迭代格式:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - [\mathbf{J}(\mathbf{x}_k)]^{-1} \mathbf{F}(\mathbf{x}_k)$$

实际计算中, 不求逆矩阵, 而是求解线性方程组:

1. 计算残差  $\mathbf{r}_k = -\mathbf{F}(\mathbf{x}_k)$ 。
2. 求解线性方程组  $\mathbf{J}(\mathbf{x}_k) \mathbf{s}_k = \mathbf{r}_k$  (得到修正量  $\mathbf{s}_k$ )。
3. 更新  $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{s}_k$ 。

## (1) 多元牛顿法的收敛性分析

**定理** (多元牛顿法局部二阶收敛). 设  $\mathbf{F}: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  满足以下条件:

1. 存在零点  $\mathbf{x}^* \in D$ , 即  $\mathbf{F}(\mathbf{x}^*) = \mathbf{0}$ 。
2.  $\mathbf{F}$  在  $\mathbf{x}^*$  的邻域内二阶连续可微 ( $C^2$ )。
3. Jacobian 矩阵  $\mathbf{J}(\mathbf{x}^*)$  非奇异 (即  $\det(\mathbf{J}(\mathbf{x}^*)) \neq 0$ ,  $\mathbf{x}^*$  为单根)。

则存在  $\mathbf{x}^*$  的邻域  $S$ , 使得对任意初始点  $\mathbf{x}_0 \in S$ , 牛顿迭代产生的序列  $\{\mathbf{x}_k\}$  收敛于  $\mathbf{x}^*$ , 且收敛阶至少为 2 (*Quadratic Convergence*)。即存在常数  $C > 0$ , 使得:

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq C \|\mathbf{x}_k - \mathbf{x}^*\|^2$$

说明 (证明思路). 1. 误差递推关系 设  $\mathbf{e}_k = \mathbf{x}_k - \mathbf{x}^*$  为第  $k$  步的误差。由牛顿法迭代公式:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - [\mathbf{J}(\mathbf{x}_k)]^{-1} \mathbf{F}(\mathbf{x}_k)$$

两边减去  $\mathbf{x}^*$ , 得:

$$\mathbf{e}_{k+1} = \mathbf{e}_k - [\mathbf{J}(\mathbf{x}_k)]^{-1} \mathbf{F}(\mathbf{x}_k)$$

2. *Taylor* 展开

将  $\mathbf{F}(\mathbf{x}^*)$  在  $\mathbf{x}_k$  处展开 (注意方向是  $\mathbf{x}^* = \mathbf{x}_k - \mathbf{e}_k$ ):

$$\mathbf{0} = \mathbf{F}(\mathbf{x}^*) = \mathbf{F}(\mathbf{x}_k) + \mathbf{J}(\mathbf{x}_k)(\mathbf{x}^* - \mathbf{x}_k) + O(\|\mathbf{x}^* - \mathbf{x}_k\|^2)$$

即:

$$\mathbf{F}(\mathbf{x}_k) = \mathbf{J}(\mathbf{x}_k) \mathbf{e}_k + O(\|\mathbf{e}_k\|^2)$$

或者更精确地写成积分余项形式:

$$\mathbf{F}(\mathbf{x}_k) = \mathbf{J}(\mathbf{x}_k) \mathbf{e}_k - \int_0^1 (1-t) \mathbf{F}''(\mathbf{x}_k - t \mathbf{e}_k)(\mathbf{e}_k, \mathbf{e}_k) dt$$

将  $\mathbf{F}(\mathbf{x}_k)$  的展开式代入误差递推式:

$$\begin{aligned} \mathbf{e}_{k+1} &= \mathbf{e}_k - [\mathbf{J}(\mathbf{x}_k)]^{-1} (\mathbf{J}(\mathbf{x}_k) \mathbf{e}_k + O(\|\mathbf{e}_k\|^2)) \\ &= \mathbf{e}_k - (\mathbf{e}_k + [\mathbf{J}(\mathbf{x}_k)]^{-1} O(\|\mathbf{e}_k\|^2)) \\ &= -[\mathbf{J}(\mathbf{x}_k)]^{-1} O(\|\mathbf{e}_k\|^2) \end{aligned}$$

两边取范数：

$$\|\mathbf{e}_{k+1}\| \leq \|[\mathbf{J}(\mathbf{x}_k)]^{-1}\| \cdot \|O(\|\mathbf{e}_k\|^2)\|$$

由于  $\mathbf{J}(\mathbf{x})$  在  $\mathbf{x}^*$  处连续且非奇异，存在  $\mathbf{x}^*$  的邻域  $S$  和常数  $M_1$ ，使得对任意  $\mathbf{x} \in S$ ， $\mathbf{J}(\mathbf{x})$  均可逆且  $\|[\mathbf{J}(\mathbf{x})]^{-1}\| \leq M_1$ （逆算子的一致有界性）。又因为  $\mathbf{F}$  二阶可微，二阶导数在闭邻域上有界，即存在常数  $M_2$  使得 Taylor 展开的余项满足  $\|O(\|\mathbf{e}_k\|^2)\| \leq M_2 \|\mathbf{e}_k\|^2$ 。于是：

$$\|\mathbf{e}_{k+1}\| \leq M_1 M_2 \|\mathbf{e}_k\|^2 = C \|\mathbf{e}_k\|^2$$

其中  $C = M_1 M_2$ 。

**结论：**这就证明了误差是按照平方速度衰减的。只要初值  $\mathbf{x}_0$  足够接近  $\mathbf{x}^*$ ，使得  $C \|\mathbf{e}_0\| < 1$ ，迭代就会收敛。

## (2) 拟牛顿法

为了避免计算昂贵的 Jacobian 矩阵及其逆，拟牛顿法使用近似矩阵  $B_k \approx \mathbf{J}(\mathbf{x}_k)$  并满足拟牛顿方程：

$$B_{k+1}(\mathbf{x}_{k+1} - \mathbf{x}_k) = \mathbf{F}(\mathbf{x}_{k+1}) - \mathbf{F}(\mathbf{x}_k)$$

**Broyden 秩 1 方法：**最常用的秩 1 更新公式：

$$B_{k+1} = B_k + \frac{(\mathbf{y}_k - B_k \mathbf{s}_k) \mathbf{s}_k^T}{\mathbf{s}_k^T \mathbf{s}_k}$$

其中  $\mathbf{s}_k = \mathbf{x}_{k+1} - \mathbf{x}_k$ ,  $\mathbf{y}_k = \mathbf{F}(\mathbf{x}_{k+1}) - \mathbf{F}(\mathbf{x}_k)$ 。

## (3) 往年题：牛顿法计算特征值与归一化特征向量

### a. 特征值问题的非线性方程组转化

对于  $n$  阶实矩阵  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ，特征值-特征向量问题  $\mathbf{Ax} = \lambda \mathbf{x} (\mathbf{x} \neq \mathbf{0})$  结合归一化约束  $\mathbf{x}^T \mathbf{x} = 1$ ，可转化为求解非线性方程组  $\mathbf{F}(\mathbf{z}) = \mathbf{0}$ ，其中  $\mathbf{z} = (\mathbf{x}^T, \lambda)^T \in \mathbb{R}^{n+1}$ ，且：

$$\mathbf{F}(\mathbf{z}) = \mathbf{F}(\mathbf{x}, \lambda) = \begin{pmatrix} (\mathbf{A} - \lambda \mathbf{I}) \mathbf{x} \\ \frac{1}{2} (\mathbf{x}^T \mathbf{x} - 1) \end{pmatrix}$$

### b. Jacobian 矩阵构造

$\mathbf{F}(\mathbf{z})$  的  $(n+1) \times (n+1)$  阶 Jacobian 矩阵为：

$$\mathbf{J}(\mathbf{z}) = \begin{pmatrix} \mathbf{A} - \lambda \mathbf{I} & -\mathbf{x} \\ \mathbf{x}^T & 0 \end{pmatrix}$$

### c. 牛顿迭代格式

待填写：() 以下构造过程就是一样的了，抄上面方程组牛顿迭代格式即可

# 六 常微分方程初值问题数值解法

## 总结

- Lipschitz 条件
- 零稳定性
- 绝对稳定性
- 局部截断误差
- 相容性与相容阶
- Taylor 展开法求局部截断误差的重要结论

## 单步法

- 一般形式
- 求解方法
  - Euler 折线法 (Euler 公式)
  - 向后 Euler 法
  - 中点格式及其显式化
  - 梯形格式及其显式化
- 一步误差
- 局部截断误差
- 相容性与相容阶

## (一) 通用理论

### 1. 问题模型

一阶初值问题提法:

$$\begin{cases} y' = f(t, y), & t \in [a, b], \quad y \in \mathbf{R}^d \\ y(a) = y_0 \end{cases}$$

半线性高阶常微分方程初值问题提法:

$$\begin{aligned} y^{(n)} + F(t, y, y', \dots, y^{(n-1)}) &= 0, \quad t \in [a, b] \\ y(a) = y_0, \quad y'(a) = y_0^{(1)}, \dots, \quad y^{(n-1)}(a) &= y_0^{(n-1)} \end{aligned}$$

从一阶初值问题到高阶初值问题的转化: 考虑导数向量

$$Y = [y, y', \dots, y^{(n-1)}]^T$$

则该导数向量满足矩阵关系:

$$Y' = \begin{bmatrix} 0 & I & 0 & \dots & 0 \\ 0 & 0 & I & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & I \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix} Y + \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ -F(t, y, y', \dots, y^{(n-1)}) \end{bmatrix}$$

### 2. 数值解法前提: Lipschitz 条件

对比: 线性方程组求解的压缩映像原理; 但这里不要求  $L < 1$ 。

**定义 (Lipschitz 条件).** 设函数  $f : [a, b] \times \mathbf{R}^d \rightarrow \mathbf{R}^d$ , 若存在常数  $L > 0$ , 使得对任意  $t \in [a, b]$  及任意  $y_1, y_2 \in \mathbf{R}^d$ , 都有

$$\|f(t, y_1) - f(t, y_2)\| \leq L \|y_1 - y_2\|$$

则称  $f$  满足 Lipschitz 条件,  $L$  为 Lipschitz 常数。

当方程满足 Lipschitz 条件时, 初值问题存在唯一解  $y = y(t) \in C^1[a, b]$

### 3. 基本概念

#### (1) 单步法与多步法的一般形式

- 单步法: 从  $n$  计算出  $n+1$

$$y_{n+1} = y_n + h_{n+1} \phi(t_n, t_{n+1}; y_n, y_{n+1}; f)$$

- 多步法：从  $(n, n+k-1)$  共  $k$  个点计算出  $n+k$

$$\sum_{i=0}^k \alpha_i^{(n+k)} y_{n+i} = h_{n+k} \phi(t_n, t_{n+1}, \dots, t_{n+k}; y_n, y_{n+1}, \dots, y_{n+k}; f)$$

式中， $\alpha_k^{n+k} = 1$

对于单步法和多步法，增量函数  $\phi$  均需要满足以下三条性质：

- 连续性： $\phi$  在所有的（除了函数  $f$  外）的变量上连续
- 零值性：当  $f = 0$  时， $\phi(\cdot; f) = 0$
- Lipschitz 条件：

– 对于单步法：

$$\|\phi(t, \tau; u_1, v_1; f) - \phi(t, \tau; u_2, v_2; f)\| \leq L_f (\|u_1 - u_2\| + \|v_1 - v_2\|)$$

– 对于多步法，需要满足一致 Lipschitz 条件：

$$\begin{aligned} & \|\phi(t_n, t_{n+1}, \dots, t_{n+k}; u_0, u_1, \dots, u_k; f) - \phi(t_n, t_{n+1}, \dots, t_{n+k}; v_0, v_1, \dots, v_k; f)\| \\ & \leq L_f \sum_{i=0}^k \|u_i - v_i\| \end{aligned}$$

## (2) 特征多项式

对于多步法，定义其特征多项式为

$$\rho(\xi) = \sum_{i=0}^k \alpha_i^{(n+k)} \xi^i$$

## (3) 零稳定性

### a. 单步法

零稳定性考虑的是方程右端函数  $f(t, y)$  受到一个小扰动变为  $\tilde{f}(t, y)$  后，数值解的变化情况。相应地，初始条件也可能被扰动，从  $y_0$  变为  $\tilde{y}_0$ .

若当扰动量控制在一个范围内时，数值解的偏差也能被控制在一个与扰动量成正比的范围内，则称该数值方法是零稳定的。

**定义 (零稳定性).** 若存在正常数  $K^*$  和  $\eta^*$ ，使得对任意  $\epsilon \leq \eta^*$ ，当对原方程施加一个扰动

$$\begin{aligned} \|y_0 - \tilde{y}_0\| &\leq \epsilon \\ \|f(t, y) - \tilde{f}(t, y)\| &\leq \epsilon, \quad t \in [a, b], y \in \mathbf{R}^d \end{aligned}$$

得到新的方程<sup>1</sup>

$$\begin{cases} v' = \tilde{f}(t, v), & t \in [a, b] \\ v(a) = y_0 \end{cases}$$

若该方程有解存在，且满足

$$\|y(t) - v(t)\| \leq K^* \epsilon, \quad t \in [a, b]$$

则称该数值方法是零稳定的。

### b. 多步法

若存在正常数  $C$  和  $h_0$ , 当  $0 < h < h_0$  时, 多步法的任意两解  $u_n$  和  $v_n$  满足不等式

$$\max |u_n - v_n| \leq C \max_{j=0,1,\dots,k-1} |u_j - v_j|$$

等价形式:

考虑  $k$  点数值向量

$$Y_n = [y_n, y_{n+1}, \dots, y_{n+k-1}]^T$$

则多步法可表示为

$$Y_{n+1} = \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \\ -\alpha_0^{(n+k)} & -\alpha_1^{(n+k)} & -\alpha_2^{(n+k)} & \dots & -\alpha_{k-1}^{(n+k)} \end{bmatrix} Y_n + h_{n+k} \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ \phi(t_n, t_{n+1}, \dots, t_{n+k}; y_n, y_{n+1}, \dots, y_{n+k}; f) \end{bmatrix}$$

## (4) 一步误差与局部截断误差

单步法:

- 一步误差: 考虑已知一点的数值解  $(t_n, y_n)$ , 设  $y = y(t; t_n, y_n)$  是点  $(t_n, y_n)$  所在的积分曲线 (即该数值解点对应的一条精确解), 用单步法得到  $y_{n+1}$  处的数值解, 则一步误差定义为由点  $(t_n, y_n)$  所确定的  $t_{n+1}$  处所确定的“一步精确值”与数值结果导出的“一步近似值”之间的差值:

$$\tilde{R}_{n+1} = y(t_{n+1}; t_n, y_n) - y_{n+1} = y(t_{n+1}; t_n, y_n) - y_n - h_{n+1} \phi(t_n, t_{n+1}; y_n, y_{n+1}; f)$$

- 局部截断误差: 在  $t_{n+1}$  处, 假设前一步的数值解是准确的, 即  $y_n = y(t_n)$ , 则在  $t_{n+1}$  的局部截断误差定义为

$$R_{n+1} = y(t_{n+1}) - y_{n+1} = y(t_{n+1}) - y(t_n) - h_{n+1} \phi(t_n, t_{n+1}; y(t_n), y(t_{n+1}); f)$$

<sup>1</sup>老爹在这里的符号规范仿佛是精神分裂了一样,  $\tilde{y}_0$  表示准确值,  $y_0$  表示扰动值;  $f$  则反过来,  $\tilde{f}$  表示扰动后的函数,  $f$  表示准确的函数。

二者满足关系

$$(1 - h_{n+1}L)|\tilde{R}_{n+1}| \leq |R_{n+1}| \leq (1 + h_{n+1}L)|\tilde{R}_{n+1}|$$

多步法：

$$R_{n+k} = * \sum_{i=0}^k \alpha_i^{(n+k)} y(t_{n+i}) - h_{n+k} \phi(t_n, t_{n+1}, \dots, t_{n+k}; y(t_n), y(t_{n+1}), \dots, y(t_{n+k}); f)$$

## (5) 相容性

单步法：对任意积分曲线  $y = y(t)$ , 若

$$R(h) = y(t+h) - y(t) - h\phi(t, t+h; y(t), y(t+h); f) = o(h)$$

则称该方法是相容的。此外，若

$$\|R(h)\| \leq Mh^{p+1}$$

则称该方法至少 p 阶相容。

若局部截断误差可以展开为

$$R(h) = \psi(t; y)h^{p+1} + O(h^{p+2})$$

则称  $\psi(t; y)h^{p+1}$  为主局部截断误差。

此外，对于隐式方法，局部截断误差和一步误差有相同的主项。

多步法：

若局部截断误差  $R^{n+k} = o(h_{n+k})$ , 则称多步法是相容的。此外，若

$$R^{n+k} = O(h_{n+k}^{p+1})$$

则称其为 p 阶相容，主项称为主局部截断误差。

## (6) 收敛性与收敛阶

定义 (单步法的收敛性). 若单步法满足

$$\sup_n \|y(t_n) - y_n\| \rightarrow 0, \quad h \rightarrow 0^+, \quad e_0 \rightarrow 0$$

定义 (多步法的收敛性). 若对任意初值

$$y_r = s_r, \quad r = 0, 1, \dots, k-1$$

满足

$$s_r(h) \rightarrow \tilde{y}_0, \quad h \rightarrow 0^+$$

时总有

$$\sup_n \|y_n - y(t_n)\| \rightarrow 0, \quad h \rightarrow 0^+$$

则称多步法是收敛的。

### (7) 稳定性

若对于给定的时间演化区间  $(0, T]$ , 存在常数  $h_0$  和  $K$ , 使得当

$$0 < h \leq h_0$$

时, 任意两数值解满足

$$\max_{t_n \leq T} |y_n - z_n| \leq K |y_0 - z_0|$$

则称单步法是零稳定的。

单步法

$$y_{n+1} = y_n + h_{n+1} \phi(t_n, t_{n+1}; y_n, y_{n+1}; f)$$

## (二) 相关定理

### 1. 单步法相容的充要条件

$$f(t, y(t)) = \phi(t, t; y(t), y(t); f)$$

### 2. 整体误差估计

记第  $n$  步的误差为  $e_n = y(t_n) - y_n$ 。若函数  $\phi$  满足前面条件, 则存在  $h_0 > 0$  和  $C > 0$ , 使得当  $0 < h < h_0$  时, 有

$$\|e_n\| \leq C \left( \|e_0\| + \sup_m \frac{\|R_m\|}{h_m} \right)$$

若方法是  $p$  阶相容的, 则

$$\|e_n\| = O(\|e_0\| + h^p)$$

### 3. 收敛性定理

#### (1) 单步法

若方法是相容的，且

$$y_0 \rightarrow \tilde{y}_0, \quad h \rightarrow 0^+$$

则方法是收敛的。

同时，若  $y_0 - \tilde{y}_0 = O(h^p)$  且方法是 p 阶相容的，则收敛阶为 p.

#### (2) 多步法

**定理** (多步法收敛的充分条件). 多步法收敛的充分条件是特征多项式在 1 处的值为 0，即  $\rho(1) = 0$

**定理** (多步法收敛性定理). 多步法若强稳定且相容，则多步法收敛。进一步，若方法 p 阶相容，并且初值的近似误差也是 p 阶的，则

$$\sup_n \|y_n - y(t_n)\| = O(h^p)$$

### 4. 稳定性定理

#### (1) 单步法

单步法是零稳定的。

(虽然 PPT 上就这么一句话，但 Gemini 说其实这要求增量函数  $\phi$  和方程右侧函数  $f$  都满足连续性和 Lipschitz 条件。)

#### (2) 多步法

##### a. 多步法的强稳定性条件

若存在某个范数，使得多步法系数矩阵的范数不大于 1，则称多步法满足强稳定性条件。

##### b. 多步法的强稳定性与零稳定性

强稳定性可以确保多步法是零稳定的。

##### c. 多步法强稳定性的一个充分条件

若对于等步长特征多项式

$$\rho(\xi) = \sum_{i=0}^k \alpha_i \xi^i$$

除了 1 这一个单根外，所有根的模值都小于 1，则存在一个常数  $\kappa > 1$ ，使得当  $h_m/h_n \leq \kappa$  时，多步法是强稳定的。

### (三) 具体求解方法

区间划分规定：

$$a = t_0 < t_1 < \dots < t_N = b$$

记号：

区间长度：

$$h_i = t_i - t_{i-1}, \quad i = 1, 2, \dots, N$$

注意，此时的区间长度为这个下标减去前一个下标，与数值微积分一节不同。

此外，记

$$h = \max_i h_i$$

准确值与近似值：

- 准确值记为  $y(t_i)$
- 近似值记为  $y_i$

## 1. 单步法

基本方法：通过各种不同数值表达式得到前一时刻和后一时刻的函数值之间的关系，进行时间步进一般形式：

$$y_{n+1} = y_n + h_{n+1}\phi(t_n, t_{n+1}; y_n, y_{n+1}; f)$$

### (1) Euler 方法

#### a. 欧拉折线法：显格式

由每一点的向前差分近似导数得到。

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

- 一阶相容性
- 主局部截断误差  $\frac{1}{2}h^2y''(t)$

### b. 向前欧拉：显格式

将问题转化为积分形式，用左矩形公式近似积分，得到 Euler 公式：

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

### c. 向后欧拉：隐格式

若用右矩形求积公式近似积分，则得到向后欧拉方法：

$$y_{n+1} = y_n + h_{n+1}f(t_{n+1}, y_{n+1})$$

## (2) 数值积分导出的其他方法

### a. 中点格式：可以显式化的隐格式

用中点求积公式近似积分，得到中点格式：

$$y_{n+1} = y_n + h_{n+1}f(t_{n+1/2}, y_{n+1/2})$$

若中点处  $y_n$  取一阶近似，即

$$\begin{aligned} t_{n+1/2} &= \frac{t_n + t_{n+1}}{2} \\ y_{n+1/2} &= y_n + \frac{h_{n+1}}{2}f(t_n, y_n) \end{aligned}$$

即为显式中点方法。

### b. 梯形公式：隐格式

用梯形公式近似积分，得到梯形格式：

$$y_{n+1} = y_n + \frac{h_{n+1}}{2}[f(t_n, y_n) + f(t_{n+1}, y_{n+1})]$$

- 二阶相容性
- 主局部截断误差  $-\frac{1}{12}h^3y^{(3)}(t)$

### c. 显式梯形公式

同样用一阶近似代入梯形公式，得到显式梯形公式：

$$\begin{aligned} y_{n+1} &= y_n + \frac{h_{n+1}}{2}[f(t_n, y_n) + f(t_{n+1}, y_{n+1}^*)] \\ y_{n+1}^* &= y_n + h_{n+1}f(t_n, y_n) \end{aligned}$$

### (3) 龙格 - 库塔 (RK) 方法

明确了必考显式龙格库塔方法的推导

思想：用  $s$  个点的斜率的线性组合来近似平均斜率，而每一个点的确定又都与其他点相关。

一般形式：

$$y_{n+1} = y_n + h_{n+1} \sum_{i=1}^s b_i k_i$$

式中，

$$k_i = f(t_n + c_i h_{n+1}, y_n + h_{n+1} \underbrace{\sum_{j=1}^s a_{ij} k_j}_{\substack{\text{用 } s \text{ 个点得到的} \\ \text{近似平均斜率}}})$$

且满足

$$\sum_{j=1}^s a_{ij} = c_i$$

一般形式下的参数或者约束：

$$\begin{aligned} c_k &= \sum_{j=1}^s a_{kj} \quad (\text{行和}) \\ \sum_{i=1}^s b_i &= 1 \quad (\text{权重参数}) \end{aligned}$$

这种一般形式的等价形式为

$$\begin{aligned} y_{n+1} &= y_n + h_{n+1} \sum_{i=1}^s b_i f(t_n + c_i h_{n+1}, Y_i) \\ Y_i &= y_n + h_{n+1} \sum_{j=1}^s a_{ij} f(t_n + c_j h_{n+1}, Y_j) \end{aligned}$$

可以理解为引入一个  $y_n$  的代数函数  $Y_i$  作为中间变量。

#### a. 显式与隐式 RK 方法

若矩阵  $A = (a_{ij})$  是严格下三角矩阵，则称该 RK 方法为显式 RK 方法，否则为隐式 RK 方法。

#### b. 显式 2 阶 RK 方法

参数组合：

$$b_1 + b_2 = 1$$

$$b_2 c_2 = \frac{1}{2}$$

### c. 3 阶 RK 方法

$$b_1 + b_2 + b_3 = 1$$

$$b_2 c_2 + b_3 c_3 = \frac{1}{2}$$

$$b_2 c_2^2 + b_3 c_3^2 = \frac{1}{3}$$

$$b_3 a_{32} c_2 = \frac{1}{6}$$

## 2. 多步法

### (1) Adams 方法

a. 显式 Adams 公式

b. 隐式 Adams 公式

### (2) Nyström 方法

### (3) Milne-Simpson 方法



# 七 线性代数方程组数值解法

## (一) 线性代数方程组直接解法 (Direct Methods)

### 1. 基本概念

求解  $Ax = b$ , 其中  $A \in \mathbb{R}^{n \times n}$  非奇异。

#### (1) Cramer 法则 (Cramer's Rule)

即直接求逆

定理 (Cramer 法则). 若  $\det(A) \neq 0$ , 则  $x_i = \det(A_i)/\det(A)$ 。

备注. 理论重要但计算不可行 ( $O((n+1)!)$ ), 仅用于极小规模或理论证明。

### 2. Gauss 消去法 (Gauss Elimination)

#### (1) 前向消去与回代

求解  $Ax = b$  分为两个阶段:

1. 消元过程 (Forward Elimination): 将增广矩阵  $[A|b]$  变换为上三角形式  $[U|y]$ 。对于第  $k$  步 ( $k = 1, 2, \dots, n-1$ ):

- 计算乘子: 设主元  $a_{kk}^{(k)} \neq 0$ , 对所有  $i = k+1, \dots, n$ :

$$m_{ik} = \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}}$$

- 行变换: 将第  $i$  行减去第  $k$  行的  $m_{ik}$  倍, 消去  $a_{ik}^{(k)}$ :

$$a_{ij}^{(k+1)} = a_{ij}^{(k)} - m_{ik}a_{kj}^{(k)}, \quad j = k+1, \dots, n$$

$$b_i^{(k+1)} = b_i^{(k)} - m_{ik}b_k^{(k)}$$

经过  $n-1$  步后, 得到上三角矩阵  $U = A^{(n)}$  和变换后的右端项  $y = b^{(n)}$ 。

2. 回代过程 (Backward Substitution): 求解上三角方程组  $Ux = y$ 。

$$\begin{cases} u_{11}x_1 + u_{12}x_2 + \cdots + u_{1n}x_n = y_1 \\ \dots \\ u_{nn}x_n = y_n \end{cases}$$

计算公式为:

$$x_n = \frac{y_n}{u_{nn}}$$

$$x_i = \frac{1}{u_{ii}} \left( y_i - \sum_{j=i+1}^n u_{ij}x_j \right), \quad i = n-1, \dots, 1$$

## (2) 选主元策略 (Pivoting)

若主元  $|a_{kk}^{(k)}|$  过小, 乘子  $m_{ik}$  会很大, 导致舍入误差剧烈放大。

- 列主元消去法: 交换行, 使  $\max_{i \geq k} |a_{ik}^{(k)}|$  所在行成为当前主行。
- 全主元消去法: 交换行和列, 选全子阵最大值。
- 数值稳定性: 选主元是保证 Gauss 消去法数值稳定的必要条件。

## 3. 矩阵三角分解 (Matrix Factorization)

### (1) LU 分解

定理 (存在唯一性). 若  $A$  的所有顺序主子式非零, 则  $A$  可唯一分解为  $A = LU$ , 其中  $L$  为单位下三角阵,  $U$  为上三角阵。

说明. 消去步的矩阵表示

在 Gauss 消去的第  $k$  步, 我们通过行变换将第  $k$  列对角线以下的元素消为 0。这一步等价于左乘一个初等下三角矩阵 (Atomic Lower Triangular Matrix)  $L_k$ :

$$L_k = I - \mathbf{m}_k \mathbf{e}_k^T = \begin{bmatrix} 1 & & & & & \\ & \ddots & & & & \\ & & 1 & & & \\ & & & -m_{k+1,k} & 1 & \\ & & & \vdots & & \ddots \\ & & & -m_{n,k} & & 1 \end{bmatrix}$$

其中  $\mathbf{m}_k = [0, \dots, 0, m_{k+1,k}, \dots, m_{n,k}]^T$  是由第  $k$  步的乘子构成的向量,  $\mathbf{e}_k$  是单位向量。

整个消去过程可以写成:

$$L_{n-1} L_{n-2} \dots L_1 A = U$$

**说明.** 逆矩阵的性质

$L_k$  的逆矩阵  $L_k^{-1}$  非常简单, 只需将乘子的符号取反:

$$L_k^{-1} = (I - \mathbf{m}_k \mathbf{e}_k^T)^{-1} = I + \mathbf{m}_k \mathbf{e}_k^T$$

这是因为  $(I - \mathbf{m}_k \mathbf{e}_k^T)(I + \mathbf{m}_k \mathbf{e}_k^T) = I - \mathbf{m}_k (\mathbf{e}_k^T \mathbf{m}_k) \mathbf{e}_k^T = I$  (注意  $\mathbf{e}_k^T \mathbf{m}_k = 0$ , 因为  $\mathbf{m}_k$  前  $k$  个元素为 0)。

**说明.** 构造  $L$  矩阵

由  $L_{n-1} \dots L_1 A = U$ , 我们有:

$$A = (L_{n-1} \dots L_1)^{-1} U = L_1^{-1} L_2^{-1} \dots L_{n-1}^{-1} U$$

令  $L = L_1^{-1} L_2^{-1} \dots L_{n-1}^{-1}$ 。神奇的是, 计算这个乘积不需要复杂的矩阵运算。由于  $L_k^{-1}$  的非零乘子位于第  $k$  列,  $L_j^{-1}$  的位于第  $j$  列 ( $j > k$ ), 它们的乘积只是简单地将各列的乘子填充到单位矩阵对应的位置, 互不干扰:

$$L = \left( I + \sum_{k=1}^{n-1} \mathbf{m}_k \mathbf{e}_k^T \right) = \begin{bmatrix} 1 & & & & \\ m_{21} & 1 & & & \\ m_{31} & m_{32} & 1 & & \\ \vdots & \vdots & \ddots & \ddots & \\ m_{n1} & m_{n2} & \dots & m_{n,n-1} & 1 \end{bmatrix}$$

结论:  $L$  矩阵的下三角部分正是 *Gauss* 消去过程中计算出的所有乘子  $m_{ik}$ 。

### a. LU 分解的应用形式

一旦得到  $A = LU$ , 求解  $Ax = b$  变为两步:

1. 解下三角方程组  $Ly = b$  (前代)。
2. 解上三角方程组  $Ux = y$  (回代)。

### (2) Cholesky 分解 ( $LL^T$ )

针对对称正定矩阵 (SPD) 的高效分解。

矩阵  $A$  存在  $LL^T$  分解 (其中  $l_{ii} > 0$ )  $\iff A$  是对称正定的。

### (3) 附: 对称正定矩阵 (SPD) 的判定方法

1. 定义法:  $A$  是对称矩阵 ( $A^T = A$ ), 且对任意非零向量  $x \neq 0$ , 都有二次型  $x^T Ax > 0$ 。
2. 特征值判别法:  $A$  是对称矩阵, 且  $A$  的所有特征值  $\lambda_i$  均严格大于零 ( $\lambda_i > 0, \forall i$ )。
3. 顺序主子式判别法 (Sylvester 准则):  $A$  的所有顺序主子式  $\det(A_k)$  均严格大于零 ( $k = 1, \dots, n$ )。

4. **Cholesky 分解存在性:** 若 Cholesky 分解算法能够顺利执行到底, 且所有对角元  $l_{ii}$  均为实数且非零 (即开方过程未遇到负数或零), 则  $A$  必为对称正定矩阵。这也是计算机程序中判断 SPD 的最实用方法。

**定理.** 若  $A$  对称正定, 则存在唯一的对角元  $> 0$  的下三角阵  $L$  使得  $A = LL^T$ 。此外, Cholesky 分解是数值稳定的, 无需选主元 (因为  $|l_{ij}|^2 \leq \sum l_{ik}^2 = a_{ii}$ , 元素不会增长)。

- **算法推导:** 比较  $A = LL^T$  两边的元素:

$$a_{ij} = \sum_{k=1}^n l_{ik}(L^T)_{kj} = \sum_{k=1}^{\min(i,j)} l_{ik}l_{jk}$$

- **计算公式** (按列计算): 对于  $j = 1, \dots, n$ :

1. 对角元:

$$l_{jj} = \sqrt{a_{jj} - \sum_{k=1}^{j-1} l_{jk}^2}$$

2. 非对角元 ( $i > j$ ):

$$l_{ij} = \frac{1}{l_{jj}} \left( a_{ij} - \sum_{k=1}^{j-1} l_{ik}l_{jk} \right)$$

- **改进:**  $LDL^T$  分解 (针对非正定对称阵): 如果  $A$  仅仅是对称的 (可能是未定或不定的), 但其顺序主子式非零, 可以进行  $A = LDL^T$  分解, 其中  $L$  是单位下三角,  $D$  是对角阵 (对角元可正可负)。

$$d_j = a_{jj} - \sum_{k=1}^{j-1} d_k l_{jk}^2, \quad l_{ij} = \frac{1}{d_j} \left( a_{ij} - \sum_{k=1}^{j-1} d_k l_{ik}l_{jk} \right)$$

这种方法避免了开方, 适用于更广泛的对称矩阵, 但若遇到  $d_j \approx 0$  仍需选主元 ( $PAP^T = LDL^T$ )。

#### (4) 追赶法 (Thomas Algorithm)

针对三对角矩阵的  $LU$  分解简化版 (待补充)。

### 4. 误差分析与条件数

#### (1) 条件数 (Condition Number)

衡量方程组解对数据扰动的敏感程度。

$$\text{cond}(A) = \|A\| \cdot \|A^{-1}\|$$

**性质.** 1.  $\text{cond}(A) \geq 1$ 。2. 若  $A$ 酉相似于  $B$ , 则  $\text{cond}_2(A) = \text{cond}_2(B)$ 。3. 对于 SPD 矩阵,  $\text{cond}_2(A) = \frac{\lambda_{\max}}{\lambda_{\min}}$ 。

## (2) 扰动误差界定理

设  $Ax = b$ ,  $A$  非奇异。考虑  $(A + \delta A)(x + \delta x) = b + \delta b$ 。若满足条件  $\|\delta A\| \cdot \|A^{-1}\| < 1$  (保证  $A + \delta A$  非奇异), 则有如下相对误差界:

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\text{cond}(A)}{1 - \text{cond}(A) \frac{\|\delta A\|}{\|A\|}} \left( \frac{\|\delta A\|}{\|A\|} + \frac{\|\delta b\|}{\|b\|} \right)$$

证明.

$$(A + \delta A)(x + \delta x) = b + \delta b$$

展开得:

$$Ax + A\delta x + \delta Ax + \delta A\delta x = b + \delta b$$

由于  $Ax = b$ , 消去得:

$$(A + \delta A)\delta x = \delta b - \delta Ax$$

两边左乘  $(A + \delta A)^{-1}$ :

$$\delta x = (A + \delta A)^{-1}(\delta b - \delta Ax)$$

利用范数放缩取范数:

$$\|\delta x\| \leq \|(A + \delta A)^{-1}\| \cdot (\|\delta b\| + \|\delta A\| \cdot \|x\|)$$

根据前面的逆矩阵扰动定理:

$$\|(A + \delta A)^{-1}\| \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \cdot \|\delta A\|}$$

代入上式:

$$\|\delta x\| \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \cdot \|\delta A\|} (\|\delta b\| + \|\delta A\| \cdot \|x\|)$$

两边除以  $\|x\|$ :

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \cdot \|\delta A\|} \left( \frac{\|\delta b\|}{\|x\|} + \|\delta A\| \right)$$

利用  $b = Ax \implies \|b\| \leq \|A\| \cdot \|x\| \implies \frac{1}{\|x\|} \leq \frac{\|A\|}{\|b\|}$ :

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \cdot \|\delta A\|} \left( \frac{\|\delta b\| \cdot \|A\|}{\|b\|} + \|\delta A\| \right)$$

引入条件数提取  $\|A\|$ , 并利用  $\text{cond}(A) = \|A\| \cdot \|A^{-1}\|$ :

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\|A^{-1}\| \cdot \|A\|}{1 - \|A^{-1}\| \cdot \|A\| \cdot \frac{\|\delta A\|}{\|A\|}} \left( \frac{\|\delta b\|}{\|b\|} + \frac{\|\delta A\|}{\|A\|} \right)$$

整理得:

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\text{cond}(A)}{1 - \text{cond}(A) \frac{\|\delta A\|}{\|A\|}} \left( \frac{\|\delta A\|}{\|A\|} + \frac{\|\delta b\|}{\|b\|} \right)$$

证毕。

□

### (3) 后验误差估计 (Posteriori Error Estimation)

问题：如何利用计算出的近似解  $\tilde{x}$  来估计真实误差  $x - \tilde{x}$ ？

设  $\tilde{x}$  为近似解，定义残差（Residual）为：

$$r = b - A\tilde{x}$$

**定理** (基于残差的误差界). 设  $A$  非奇异，则相对误差与相对残差之间满足如下不等式：

$$\frac{1}{\text{cond}(A)} \frac{\|r\|}{\|b\|} \leq \frac{\|x - \tilde{x}\|}{\|x\|} \leq \text{cond}(A) \frac{\|r\|}{\|b\|}$$

**证明.** 1. 上界：由  $r = b - A\tilde{x} = Ax - A\tilde{x} = A(x - \tilde{x})$ ，可得：

$$x - \tilde{x} = A^{-1}r \implies \|x - \tilde{x}\| \leq \|A^{-1}\| \cdot \|r\|$$

又由  $Ax = b \implies \|b\| \leq \|A\| \cdot \|x\| \implies \frac{1}{\|x\|} \leq \frac{\|A\|}{\|b\|}$ 。两式相乘：

$$\frac{\|x - \tilde{x}\|}{\|x\|} \leq \|A^{-1}\| \cdot \|r\| \cdot \frac{\|A\|}{\|b\|} = \text{cond}(A) \frac{\|r\|}{\|b\|}$$

2. 下界：由  $r = A(x - \tilde{x}) \implies \|r\| \leq \|A\| \cdot \|x - \tilde{x}\| \implies \|x - \tilde{x}\| \geq \frac{\|r\|}{\|A\|}$ 。又由  $x = A^{-1}b \implies \|x\| \leq \|A^{-1}\| \cdot \|b\| \implies \frac{1}{\|x\|} \geq \frac{1}{\|A^{-1}\| \cdot \|b\|}$ 。两式相乘：

$$\frac{\|x - \tilde{x}\|}{\|x\|} \geq \frac{\|r\|}{\|A\| \cdot \|A^{-1}\| \cdot \|b\|} = \frac{1}{\text{cond}(A)} \frac{\|r\|}{\|b\|}$$

证毕。 □

**备注** (物理意义). 1. 残差小  $\neq$  误差小：如果  $\text{cond}(A)$  很大（矩阵病态），即使残差  $\|r\|$  很小（计算机解出的方程组几乎“满足”方程），解的误差  $\|x - \tilde{x}\|$  也可能非常大。这说明对于病态方程组，仅仅看残差是不够的。

2. **迭代改善 (Iterative Refinement):** 利用残差可以提高解的精度。计算  $r = b - A\tilde{x}$  (需用高精度计算)，解修正方程  $Ae = r$ ，则新解  $x_{\text{new}} = \tilde{x} + e$  会更精确。

### (4) 病态矩阵 (Ill-conditioned Matrix)

条件数很大的矩阵。

**例题** (Hilbert 矩阵).

$$H_{ij} = \frac{1}{i + j - 1}$$

$H_n$  高度病态。当  $n$  稍大时， $\text{cond}(H_n)$  极大，直接法解将完全失真。

### 5. 广义逆与正则化 (Generalized Inverse & Regularization)

针对奇异或极度病态方程组 (参考第二章讲义)

## (1) Moore-Penrose 广义逆

对于  $A \in \mathbb{C}^{m \times n}$  (秩  $r$ )，存在唯一的  $A^\dagger$  满足 Penrose 方程。

- **最小二乘解:** 方程组  $Ax = b$  的最小二乘解集为  $S = \{x | A^H Ax = A^H b\}$ 。
- **广义逆解:**  $x = A^\dagger b$  是  $S$  中 2-范数最小的解。
- **SVD 计算:** 若  $A = U\Sigma V^H$ ，则

$$A^\dagger = V\Sigma^\dagger U^H$$

其中  $\Sigma^\dagger = \text{diag}(\sigma_1^{-1}, \dots, \sigma_r^{-1}, 0, \dots, 0)$ 。

## (2) Tikhonov 正则化

用于处理严重的病态问题（如第一类 Fredholm 积分方程离散化）。

- **变分原理:** 最小化带罚项的泛函：

$$J_\alpha(x) = \|Ax - b\|_2^2 + \alpha\|x\|_2^2$$

其中  $\alpha > 0$  为正则化参数。

- **正规方程:**

$$(\alpha I + A^H A)x_\alpha = A^H b$$

- **效果:** 矩阵  $\alpha I + A^H A$  的条件数比  $A^H A$  改善很多（特征值整体平移  $\alpha$ ），解更稳定。
- **收敛性:** 当  $\alpha \rightarrow 0$  时， $x_\alpha \rightarrow A^\dagger b$ 。

# (二) 单步定常线性迭代解法 (Stationary Linear Iterative Methods)

## 1. 迭代法基础理论

### (1) 迭代法的基本概念

- **构造思想:** 对于大规模稀疏方程组  $Ax = b$ ，直接法计算代价过高。迭代法通过构造向量序列  $\{x^{(k)}\}$ ，使其极限为方程的解。

$$\lim_{k \rightarrow \infty} x^{(k)} = x^*$$

- 单步定常线性迭代：迭代格式为：

$$x^{(k+1)} = Bx^{(k)} + f$$

其中  $B$  称为迭代矩阵， $f$  为常向量。这种格式只依赖于前一步  $x^{(k)}$ （单步），且  $B$  和  $f$  不随  $k$  变化（定常）。

- 构造方法（矩阵分裂）：将  $A$  分裂为  $A = M - N$ ，其中  $M$  是非奇异的“近似”矩阵（易于求逆，如对角阵或三角阵）。

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

由此导出的迭代格式为：

$$Mx^{(k+1)} = Nx^{(k)} + b \implies x^{(k+1)} = M^{-1}Nx^{(k)} + M^{-1}b$$

此时迭代矩阵  $B = M^{-1}N$ ，常向量  $f = M^{-1}b$ 。

- Richardson 迭代：最简单的迭代格式，取  $M = \frac{1}{\omega}I$ 。

$$x^{(k+1)} = x^{(k)} + \omega(b - Ax^{(k)}) = (I - \omega A)x^{(k)} + \omega b$$

迭代矩阵  $B = I - \omega A$ 。收敛条件： $0 < \omega < \frac{2}{\lambda_{\max}(A)}$ （假设  $A$  为 SPD）。

- 预处理迭代：若引入预处理矩阵  $M \approx A$ ，则 Richardson 迭代变为：

$$x^{(k+1)} = x^{(k)} + \omega M^{-1}(b - Ax^{(k)})$$

这等价于对预处理后的方程  $M^{-1}Ax = M^{-1}b$  进行 Richardson 迭代。

## (2) 向量序列的收敛性

**定义** (序列收敛). 设  $\{x^{(k)}\}_{k=0}^{\infty}$  是赋范空间  $(\mathbb{R}^n, \|\cdot\|)$  中的向量序列。若存在  $x \in \mathbb{R}^n$ ，使得

$$\lim_{k \rightarrow \infty} \|x^{(k)} - x\| = 0$$

则称序列  $\{x^{(k)}\}$  收敛于  $x$ ，记为  $\lim_{k \rightarrow \infty} x^{(k)} = x$ 。

**定理** (收敛性的等价性). 1. 范数无关性：由于有限维空间  $\mathbb{R}^n$  上的所有范数都是等价的，因此向量序列的收敛性与具体选择的范数无关。如果一个序列在某种范数下收敛，它在任意范数下都收敛。

2. 分量收敛性：向量序列  $\{x^{(k)}\}$  收敛于  $x$  的充要条件是其每一个分量序列都收敛于  $x$  的对应分量。

$$\lim_{k \rightarrow \infty} x^{(k)} = x \iff \lim_{k \rightarrow \infty} x_i^{(k)} = x_i, \quad \forall i = 1, \dots, n$$

### (3) 迭代误差与收敛速度

设精确解为  $x^*$ , 第  $k$  步误差为  $e^{(k)} = x^{(k)} - x^*$ 。

- 误差传播方程:

$$e^{(k+1)} = x^{(k+1)} - x^* = (Bx^{(k)} + f) - (Bx^* + f) = B(x^{(k)} - x^*) = Be^{(k)}$$

递推可得:

$$e^{(k)} = B^k e^{(0)}$$

- **Jordan 标准型分析与收敛充要条件:** 将  $B$  转化为 Jordan 标准型  $B = PJP^{-1}$ , 则  $B^k = PJ^kP^{-1}$ 。对于 Jordan 块  $J_i(\lambda_i)$ :

$$J_i^k = \begin{bmatrix} \lambda_i^k & {}_{(1)}^k \lambda_i^{k-1} & \dots & \\ & \lambda_i^k & {}_{(1)}^k \lambda_i^{k-1} & \\ & & \ddots & \end{bmatrix}$$

要使  $B^k \rightarrow 0$  (即  $e^{(k)} \rightarrow 0$  对任意初值成立), 充要条件是所有特征值的模  $|\lambda_i| < 1$ 。即 **谱半径条件:**  $\rho(B) < 1$ 。

- **平均收敛率 (Average Rate of Convergence):** 考查前  $k$  步误差范数的平均衰减:

$$\frac{\|e^{(k)}\|}{\|e^{(0)}\|} \leq \|B^k\| \approx \rho(B)^k$$

定义  $k$  步平均收敛率为:

$$R_k(B) = -\frac{1}{k} \ln \|B^k\|$$

- **渐进收敛率 (Asymptotic Rate of Convergence):** 当  $k \rightarrow \infty$  时, 利用 Gelfand 半径公式  $\lim_{k \rightarrow \infty} \|B^k\|^{1/k} = \rho(B)$ , 可得:

$$R(B) = \lim_{k \rightarrow \infty} R_k(B) = -\ln \rho(B)$$

**意义:**  $R(B)$  越大 (即  $\rho(B)$  越小), 收敛越快。若要使误差减小  $10^{-m}$  倍 (增加  $m$  位有效数字), 所需迭代次数  $k$  约为:

$$k \approx \frac{m \ln 10}{R(B)}$$

## 2. 经典迭代方法

将  $A$  分解为  $A = D - L - U$ :

- $D$ : 对角部分。
- $-L$ : 严格下三角部分。
- $-U$ : 严格上三角部分。

### (1) Jacobi 迭代

- 分裂:  $M = D, N = L + U$ 。
- 公式:  $Dx^{(k+1)} = (L + U)x^{(k)} + b$ 。
- 分量形式:

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

- 迭代矩阵:  $B_J = D^{-1}(L + U)$ 。

### (2) Gauss-Seidel 迭代

- 思想: 计算  $x_i^{(k+1)}$  时, 利用已更新的  $x_1^{(k+1)}, \dots, x_{i-1}^{(k+1)}$ 。
- 分裂:  $M = D - L, N = U$ 。
- 公式:  $(D - L)x^{(k+1)} = Ux^{(k)} + b$ 。
- 分量形式:

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

- 迭代矩阵:  $B_{GS} = (D - L)^{-1}U$ 。
- 收敛性: 若  $A$  为 SPD 矩阵, GS 法必收敛。

### (3) 超松弛迭代 (SOR 法)

- 思想: 对 GS 迭代进行加权平均 (外推), 加速收敛。

$$x^{(k+1)} = (1 - \omega)x^{(k)} + \omega x_{GS}^{(k+1)}$$

其中  $\omega$  为松弛因子。

- 分裂:  $M = \frac{1}{\omega}(D - \omega L)$ 。
- 迭代矩阵:

$$B_{SOR} = (D - \omega L)^{-1}[(1 - \omega)D + \omega U]$$

- 收敛性:

1. 必要条件:  $0 < \omega < 2$  (Kahan 定理)。
2. 若  $A$  为 SPD, 则  $0 < \omega < 2 \iff SOR$  收敛。

- **最优松弛因子:** 对于具有 Property A 的矩阵 (如三对角阵):

$$\omega_{opt} = \frac{2}{1 + \sqrt{1 - \rho(B_J)^2}}$$

此时  $\rho(B_{SOR}) = \omega_{opt} - 1$ 。

### (三) 非定常迭代法 (Krylov Subspace Methods)

#### 1. 变分原理与 Ritz 方法基础

##### (1) 算子方程与变分问题的等价性

考虑线性方程组  $Ax = b$ , 其中  $A \in \mathbb{R}^{n \times n}$  是 \*\* 对称正定 (SPD)\*\* 矩阵。

**定理 (变分原理).** 求解  $Ax = b$  等价于寻找二次泛函  $\phi(x)$  的极小值点:

$$\phi(x) = \frac{1}{2}(Ax, x) - (b, x) = \frac{1}{2}x^T Ax - b^T x$$

即  $x^* = A^{-1}b \iff x^* = \arg \min_{x \in \mathbb{R}^n} \phi(x)$ 。

证明. 计算  $\phi(x)$  的梯度:

$$\nabla \phi(x) = \frac{1}{2}(A + A^T)x - b = Ax - b \quad (\text{因 } A = A^T)$$

令梯度为零, 即  $Ax - b = 0 \implies Ax = b$ 。由于  $A$  正定, Hessian 矩阵  $\nabla^2 \phi(x) = A$  正定, 故驻点为严格极小值点。  $\square$

##### (2) Ritz 方法与子空间逼近

- **基本思想:** 在全空间  $\mathbb{R}^n$  中求极值可能太慢。Ritz 方法试图在一系列逐渐扩大的低维子空间  $V_k$  中寻找  $\phi(x)$  的近似极小点。
- **Ritz 问题:** 设  $V_k = \text{span}\{p_0, p_1, \dots, p_{k-1}\}$ 。寻找  $x_k \in x_0 + V_k$  使得

$$\phi(x_k) = \min_{x \in x_0 + V_k} \phi(x)$$

- **几何解释:**  $x_k$  是精确解  $x^*$  在子空间  $x_0 + V_k$  上的 \*\*A-正交投影\*\*。即误差  $e_k = x_k - x^*$  满足  $e_k \perp_A V_k$  ( $e_k^T A v = 0, \forall v \in V_k$ )。

### (3) 梯度与下降方向

- **残差与梯度:**  $r(x) = b - Ax = -\nabla\phi(x)$ 。残差即负梯度方向。
- **最速下降方向:**  $p_k = r_k$ 。这是局部下降最快的方向，但在椭球狭长山谷中会导致“锯齿形”振荡，收敛极慢。
- **A-共轭方向:** 为克服锯齿效应，我们希望搜索方向  $p_0, p_1, \dots$  能够使得每一步的搜索互不干扰。这就引入了 \*\*A-正交（共轭）\*\* 的概念：

$$(p_i, p_j)_A = p_i^T A p_j = 0, \quad i \neq j$$

在 A-共轭方向系下，二次型  $\phi(x)$  的等高线被“拉圆”了，从而可以快速收敛。

### (4) 一维搜索理论 (One-Dimensional Search)

在变分方法中，给定当前点  $x_k$  和下降方向  $p_k$ ，我们需要确定一个步长  $\alpha_k$ ，使得目标函数在沿该方向上达到极小。

$$x_{k+1} = x_k + \alpha_k p_k$$

- **一维搜索问题:** 寻找  $\alpha_k$  使得  $f(\alpha) = \phi(x_k + \alpha p_k)$  最小。
- **最优步长推导:** 将  $f(\alpha)$  对  $\alpha$  求导并令其为 0：

$$f'(\alpha) = \frac{d}{d\alpha} \phi(x_k + \alpha p_k) = (\nabla\phi(x_k + \alpha p_k), p_k)$$

代入  $\nabla\phi(x) = Ax - b = -r(x)$ ，得：

$$-(r(x_k + \alpha p_k), p_k) = 0$$

利用残差的线性性质  $r(x_k + \alpha p_k) = b - A(x_k + \alpha p_k) = (b - Ax_k) - \alpha Ap_k = r_k - \alpha Ap_k$ ：

$$(r_k - \alpha Ap_k, p_k) = 0 \implies (r_k, p_k) - \alpha (Ap_k, p_k) = 0$$

解得最优步长：

$$\alpha_k = \frac{(r_k, p_k)}{(Ap_k, p_k)}$$

- **结论:** 无论是在最速下降法还是共轭梯度法中，每一步的步长  $\alpha_k$  都是由上述公式确定的，它保证了新残差  $r_{k+1}$  与当前搜索方向  $p_k$  正交。

### (5) 单调性与收敛性分析

**命题:** 在给定下降方向  $p_k$  (满足  $p_k^T r_k > 0$ ) 并采用最优步长  $\alpha_k$  后，目标函数  $\phi(x)$ 、残差范数 (在特定条件下) 及误差范数 (在特定度量下) 均呈现单调下降趋势。

证明. 1. 目标函数值下降 ( $\phi(x_{k+1}) < \phi(x_k)$ ): 将  $x_{k+1} = x_k + \alpha_k p_k$  代入二次泛函展开:

$$\phi(x_{k+1}) = \phi(x_k) + \alpha_k(Ax_k - b, p_k) + \frac{1}{2}\alpha_k^2(Ap_k, p_k)$$

由于  $Ax_k - b = -r_k$ , 所以  $(Ax_k - b, p_k) = -(r_k, p_k)$ 。代入最优步长  $\alpha_k = \frac{(r_k, p_k)}{(Ap_k, p_k)}$ :

$$\begin{aligned}\phi(x_{k+1}) &= \phi(x_k) - \frac{(r_k, p_k)^2}{(Ap_k, p_k)} + \frac{1}{2} \left( \frac{(r_k, p_k)}{(Ap_k, p_k)} \right)^2 (Ap_k, p_k) \\ &= \phi(x_k) - \frac{(r_k, p_k)^2}{2(Ap_k, p_k)}\end{aligned}$$

因为  $A$  正定,  $(Ap_k, p_k) > 0$ , 且  $(r_k, p_k) \neq 0$  (否则已收敛), 故减项恒正。

$$\therefore \phi(x_{k+1}) < \phi(x_k)$$

这证明了变分函数值是严格单调下降的。

2. 误差的  $A$ -范数下降: 定义误差  $e_k = x^* - x_k$ 。注意这里  $\phi(x)$  与误差的  $A$ -范数有直接关系:

$$\phi(x) = \frac{1}{2}(A(x^* - e), x^* - e) - (b, x^* - e) = \phi(x^*) + \frac{1}{2}(Ae, e) = \phi(x^*) + \frac{1}{2}\|e\|_A^2$$

由于  $\phi(x^*) = \text{const}$  且  $\phi(x_{k+1}) < \phi(x_k)$ , 直接推出:

$$\|e_{k+1}\|_A < \|e_k\|_A$$

即误差在  $A$ -范数意义下是单调递减的。

3. 残差的正交性: 由一维搜索条件  $f'(\alpha_k) = 0$  可知:

$$(r_{k+1}, p_k) = 0$$

这意味着新残差与当前的搜索方向正交。在几何上, 这是在此方向上能达到的“最低点”的特征。□

## 2. 最速下降法 (Steepest Descent)

### (1) 算法流程与性质

- 下降方向: 取负梯度方向  $p_k = r_k = b - Ax_k$ 。

- 最优步长:  $\alpha_k = \frac{(r_k, r_k)}{(Ar_k, r_k)}$ 。

- 迭代公式:

$$x_{k+1} = x_k + \alpha_k r_k$$

- 残差递推公式:

$$r_{k+1} = r_k - \alpha_k Ar_k$$

这避免了每次迭代计算  $b - Ax_{k+1}$  的矩阵向量乘法。

**定理** (相邻残差正交性). 在最速下降法中, 相邻两步的残差相互正交, 即  $(r_{k+1}, r_k) = 0$ 。

证明. 由最优步长条件 (一维搜索性质) 可知  $(r_{k+1}, p_k) = 0$ 。而在最速下降法中,  $p_k = r_k$ 。故  $(r_{k+1}, r_k) = 0$ 。  $\square$

**备注** (锯齿现象). 由于  $r_{k+1} \perp r_k$ , 意味着搜索方向在每一步都必须拐 90 度弯。在等高线为狭长椭球 (条件数  $\kappa$  很大) 的情况下, 这会导致算法在“山谷”底部反复震荡, 前进极慢, 形成锯齿形路径。

## (2) 收敛性定理与证明 (Kantorovich 不等式)

**定理** (最速下降法收敛速度). 设  $A$  是 SPD 矩阵, 特征值为  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n > 0$ 。则最速下降法的误差满足:

$$\|e_{k+1}\|_A \leq \frac{\kappa - 1}{\kappa + 1} \|e_k\|_A$$

其中  $\kappa = \text{cond}_2(A) = \frac{\lambda_1}{\lambda_n} = \frac{\lambda_{\max}}{\lambda_{\min}}$ 。

说明. 证明思路 1. 误差递推关系: 由  $x_{k+1} = x_k + \alpha_k r_k$ , 两边减  $x^*$  得  $e_{k+1} = e_k - \alpha_k r_k$ 。注意  $r_k = b - Ax_k = A(x^* - x_k) = Ae_k$ , 代入得:

$$e_{k+1} = e_k - \alpha_k Ae_k = (I - \alpha_k A)e_k$$

2. 能量范数递推: 计算  $\|e_{k+1}\|_A^2 = (Ae_{k+1}, e_{k+1})$ :

$$\|e_{k+1}\|_A^2 = (A(e_k - \alpha_k r_k), e_k - \alpha_k r_k)$$

展开并代入最优步长  $\alpha_k = \frac{(r_k, r_k)}{(Ar_k, r_k)}$  (推导略繁琐, 直接利用函数值下降公式): 我们已知  $\phi(x_{k+1}) = \phi(x_k) - \frac{(r_k, r_k)^2}{2(Ar_k, r_k)}$ 。利用  $\phi(x) = \phi(x^*) + \frac{1}{2}\|e\|_A^2$ , 得:

$$\frac{1}{2}\|e_{k+1}\|_A^2 = \frac{1}{2}\|e_k\|_A^2 - \frac{(r_k, r_k)^2}{2(Ar_k, r_k)}$$

于是:

$$\frac{\|e_{k+1}\|_A^2}{\|e_k\|_A^2} = 1 - \frac{(r_k, r_k)^2}{(Ar_k, r_k) \cdot \|e_k\|_A^2}$$

注意  $\|e_k\|_A^2 = (Ae_k, e_k) = (r_k, A^{-1}r_k)$ , 代入上式:

$$\frac{\|e_{k+1}\|_A^2}{\|e_k\|_A^2} = 1 - \frac{(r_k, r_k)^2}{(Ar_k, r_k)(A^{-1}r_k, r_k)}$$

3. 利用 Kantorovich 不等式: 对于 SPD 矩阵  $A$ , Kantorovich 不等式指出:

$$\frac{(x, x)^2}{(Ax, x)(A^{-1}x, x)} \geq \frac{4\lambda_1\lambda_n}{(\lambda_1 + \lambda_n)^2}$$

将  $x$  替换为  $r_k$ , 代入递推式:

$$\frac{\|e_{k+1}\|_A^2}{\|e_k\|_A^2} \leq 1 - \frac{4\lambda_1\lambda_n}{(\lambda_1 + \lambda_n)^2} = \left(\frac{\lambda_1 - \lambda_n}{\lambda_1 + \lambda_n}\right)^2 = \left(\frac{\kappa - 1}{\kappa + 1}\right)^2$$

开方即得证。

### 3. 共轭梯度法 (Conjugate Gradient, CG)

#### (1) 基本思想

- **共轭方向:** 寻找一组关于  $A$  正交 (共轭) 的方向  $\{p_0, p_1, \dots\}$ , 即  $p_i^T A p_j = 0$  ( $i \neq j$ )。
- **最优性:** 第  $k$  步得到的解  $x_k$  使得  $\phi(x)$  在 Krylov 子空间  $\mathcal{K}_k(A, r_0) = \text{span}\{r_0, Ar_0, \dots, A^{k-1}r_0\}$  上达到极小。

#### (2) CG 算法流程

初始化  $x_0$ , 计算  $r_0 = b - Ax_0$ , 令  $p_0 = r_0$ 。对于  $k = 0, 1, \dots$ :

1.  $\alpha_k = \frac{r_k^T r_k}{p_k^T A p_k}$  (计算步长)
2.  $x_{k+1} = x_k + \alpha_k p_k$  (更新解)
3.  $r_{k+1} = r_k - \alpha_k A p_k$  (更新残差)
4.  $\beta_k = \frac{r_{k+1}^T r_{k+1}}{r_k^T r_k}$  (计算方向更新系数)
5.  $p_{k+1} = r_{k+1} + \beta_k p_k$  (更新搜索方向)

#### (3) 收敛性分析

- **有限步终止:** 理论上至多  $n$  步收敛到精确解 (无舍入误差时)。

- **误差估计:**

$$\frac{\|e_k\|_A}{\|e_0\|_A} \leq 2 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k$$

收敛速度远快于最速下降法, 且依赖于  $\sqrt{\text{cond}(A)}$ 。

### 4. 预处理技术 (Preconditioning)

- **目的:** 改善矩阵条件数, 加速 CG 收敛。
- **方法:** 求解等价方程  $M^{-1}Ax = M^{-1}b$ , 其中  $M \approx A$  且  $M$  易求逆。
- **预优 CG (PCG):** 在 CG 算法中引入  $M^{-1}$ , 实际上是在  $M$ -内积下进行正交化。
- **常用预条件子:** Jacobi (对角 Scaling), Incomplete Cholesky (IC)。

**待填写: ()** 从 PPT 和作业来看, 预处理 CG 的细节不是考试重点, 大概考不了... 艰深而收获小, 摆了, 建议大家多花点精力复习别的章节 (