

Numerical Linear Algebra and Algorithms
MATLAB Edition

数值线性代数与算法 (MATLAB版)



©马昌凤 柯艺芬 唐嘉 陈宝国 编著



國防工業出版社
National Defense Industry Press

责任编辑：丁福志
责任校对：苏向颖
封面设计：蒋秀芹

ding@ndip.cn



本书各章源程序可到国防工业出版社“资源下载”栏目下载，
或发邮件到896369667@qq.com索取



► 上架建议：数值代数 ◀

<http://www.ndip.cn>

ISBN 978-7-118-11320-4

9 787118 113204 >

定价：59.00 元

数值线性代数与算法

(MATLAB 版)

马昌凤 柯艺芬

编著

唐 嘉 陈宝国

国防工业出版社

· 北京 ·

内 容 简 介

本书较为系统地介绍了数值线性代数的基本理论、方法及其主要算法的 MATLAB 程序实现. 全书共分为 7 章, 内容包括矩阵代数基础、正交变换和投影方法、线性方程组的矩阵分裂迭代法、线性方程组的 Krylov 子空间迭代法、线性最小二乘问题的数值解法、解线性方程组的直接法和矩阵特征值问题的数值方法. 书中配有丰富的例题和习题, 可供学习者使用. 本书既注意保持理论分析的严谨性, 又注重计算方法的实用性, 强调算法的 MATLAB 程序在计算机上的实现.

本书内容新颖, 叙述流畅, 可作为高等学校数学与应用数学和信息与计算科学专业高年级本科生教材, 特别适用于计算数学专业研究生“数值线性代数”课程的教材或参考书, 也可供理工科其他有关专业的研究生和对数值代数与算法感兴趣的工程技术人员参考使用.

图书在版编目 (CIP) 数据

数值线性代数与算法: MATLAB 版 / 马昌凤等编著. —北京: 国防工业出版社, 2017.6
ISBN 978-7-118-11320-4

I. ①数… II. ①马… III. ①Matlab 软件—应用—线性代数数算法 IV. ①O241.6-39

中国版本图书馆 CIP 数据核字 (2017) 第 079784 号

※

国防工业出版社 出版发行

(北京市海淀区紫竹院南路 23 号 邮政编码 100048)

腾飞印务有限公司印刷

新华书店经售

*

开本 787×1092 1/16 印张 26 字数 629 千字
2017 年 6 月第 1 版第 1 次印刷 印数 1—3000 册 定价 59.00 元

(本书如有印装错误, 我社负责调换)

国防书店: (010)88540777

发行邮购: (010)88540776

发行传真: (010)88540755

发行业务: (010)88540717

前 言

数值线性代数是研究代数问题 (线性代数方程组和矩阵特征值问题) 的数值计算方法及其有关理论的一门学科. 可以说科学与工程计算领域的相当一部分应用问题最终都会归结为线性代数方程组的数值求解问题, 或矩阵的特征值与特征向量的计算问题. 数值代数是一门理论性和实际应用性都很强的学科, 并且随着计算机技术的发展, 能够进行数值计算的实际问题的规模不断增大, 相应的数值方法也在不断地改进和创新. 因此, 对于从事科学计算的大学生、研究生甚至工程技术人员来说, 系统地了解和掌握数值代数的基本理论和方法, 特别是新近发展起来且较为成熟的新型算法是非常重要的.

本书系统地介绍了数值线性代数的三大分支, 即线性代数方程组的解法、线性最小二乘法和矩阵特征值问题, 内容包括: 数值线性代数理论基础、正交变换和投影方法、线性方程组的矩阵分裂迭代法、线性方程组的 Krylov 子空间迭代法、线性最小二乘问题的数值解法、解线性方程组的直接法、矩阵特征值问题的数值方法等. 对所讨论的方法, 除了对其收敛性及计算过程的稳定性有较详尽的论述外, 还特别注重这些算法的 MATLAB 程序在计算机上的实现. 本书可作为高等学校数学与应用数学和信息与计算科学专业高年级本科生教材, 特别适用于计算数学专业研究生“数值线性代数”课程的教材或参考书, 也可供理工科其他有关专业的研究生和对数值代数与算法感兴趣的工程技术人员参考使用. 读者只需具备微积分、线性代数和 MATLAB 程序设计方面的初步知识即可顺利阅读.

本书各章节的主要算法都给出了 MATLAB 程序及相应的计算实例. 为了更好地配合教学或自学, 作者编制了与本书配套的电子课件 (PDF 格式的 PPT) 和全部算法的 MATLAB 程序, 需要的读者可到国防工业出版社网站“资源下载”栏目下载 (网址: <http://www.ndip.cn>), 或发邮件至 896369667@qq.com 索取.

由于作者水平有限, 加之时间仓促, 书中的缺点和错误在所难免, 恳请读者不吝赐教. 来信请发至 macf88@163.com 或 keyifen2017@163.com.

作 者

2017 年 3 月

目 录

第 1 章 数值线性代数理论基础	1
1.1 一些概念和记号	1
1.2 几种常用的矩阵分解	4
1.2.1 矩阵的特征分解	4
1.2.2 矩阵的 Schur 分解	6
1.2.3 矩阵的奇异值分解	12
1.2.4 矩阵的极分解和满秩分解	16
1.3 向量和矩阵的范数	19
1.3.1 向量内积与向量范数	19
1.3.2 矩阵范数与内积	21
1.4 矩阵的广义逆	28
1.5 几种特殊的矩阵类型	31
1.6 模型问题: Poisson 问题	35
习题 1	37
第 2 章 正交变换和投影方法	39
2.1 两种常用的正交变换	39
2.1.1 Householder 变换	39
2.1.2 Givens 变换	45
2.2 QR 分解	49
2.2.1 Householder 变换 QR 分解	49
2.2.2 Givens 变换 QR 分解	53
2.3 线性无关向量组的正交化	58
2.3.1 Gram-Schmidt 正交化	58
2.3.2 Householder 正交化	61
2.4 Krylov 子空间及其正交化	64
2.4.1 Krylov 子空间	64
2.4.2 Arnoldi 正交分解	66
2.4.3 Lanczos 正交分解	71
2.5 投影方法	73
2.5.1 投影算子及其性质	73
2.5.2 投影方法的基本框架	76
2.5.3 一维投影方法	80
习题 2	83

第 3 章 线性方程组的矩阵分裂迭代法	85
3.1 迭代法的一般理论	85
3.1.1 迭代法的定义与分类	85
3.1.2 收敛性与收敛速度	86
3.1.3 相容性和敏感性分析	89
3.1.4 几种常见的矩阵分裂	91
3.2 几种经典迭代法	93
3.2.1 Richardson 迭代法	93
3.2.2 Jacobi 迭代法	94
3.2.3 Gauss-Seidel (GS) 迭代法	98
3.3 松弛型迭代法	102
3.3.1 SOR 迭代法	102
3.3.2 SSOR 迭代法	106
3.3.3 AOR 迭代法	109
3.4 HSS 迭代法	111
3.4.1 HSS 和 IHSS 方法	111
3.4.2 PHSS 迭代法	119
3.5 迭代法的加速方法	124
3.5.1 外推方法	124
3.5.2 整体校正方法	126
3.5.3 基于矩阵特征值的外推方法	130
3.5.4 Chebyshev 加速方法	132
3.6 块三对角方程组的迭代解法	137
3.6.1 $PE(\alpha)$ 方法	137
3.6.2 二次 $PE(\alpha)$ 方法	141
习题 3	143
第 4 章 线性方程组的 Krylov 子空间迭代法	145
4.1 共轭梯度法	145
4.1.1 基本 CG 方法	146
4.1.2 收敛性分析	152
4.1.3 预处理 CG 方法	157
4.1.4 CGNR 方法和 CGNE 方法	159
4.2 广义极小残量法	162
4.2.1 GMRES 方法	162
4.2.2 预处理 GMRES 方法	169
4.2.3 收敛性分析	172

4.3	极小残量法	181
4.3.1	MINRES 方法	181
4.3.2	PMINRES 方法	188
4.3.3	收敛性分析	197
4.4	SYMMLQ 方法	198
4.4.1	SYMMLQ 方法	198
4.4.2	收敛性分析	203
4.5	拟极小残量法	206
4.5.1	非对称 Lanczos 方法	207
4.5.2	QMRES 方法	211
4.6	LSQR 方法	216
4.6.1	Lanczos 双对角化方法	217
4.6.2	LSQR 算法	219
4.7	广义共轭残量法	223
4.7.1	GCR 方法	224
4.7.2	GCR(m) 方法	230
4.8	投影类方法	233
4.8.1	BCG 方法	233
4.8.2	CGS 方法	238
4.8.3	BCGSTAB 方法	241
	习题 4	246
第 5 章	线性最小二乘问题的数值解法	247
5.1	线性最小二乘问题的数学性质	247
5.1.1	最小二乘解的特征及一般表示	247
5.1.2	线性 LS 的等价性问题	250
5.1.3	线性最小二乘问题的正则化	251
5.2	求解满秩最小二乘问题的数值方法	253
5.2.1	法方程方法	254
5.2.2	QR 分解方法	254
5.3	求解秩亏最小二乘问题的数值解法	256
5.3.1	列主元 QR 分解法	256
5.3.2	奇异值分解法	261
5.4	求解最小二乘问题的迭代方法	262
5.4.1	基于法方程的矩阵分裂迭代法	262
5.4.2	基于法方程的共轭梯度法	267
5.4.3	基于 KKT 方程的 SOR 类迭代法	270

5.4.4 基于 KKT 方程的 HSS 迭代法	275
习题 5	279
第 6 章 解线性方程组的直接法	281
6.1 Gauss 消去法	281
6.1.1 顺序 Gauss 消去法	281
6.1.2 列主元 Gauss 消去法	285
6.2 LU 分解法	288
6.2.1 顺序 LU 分解法	289
6.2.2 列主元 LU 分解法	291
6.2.3 不完全 LU 分解	295
6.3 对称正定方程组的直接法	298
6.3.1 Cholesky 分解法	299
6.3.2 不完全 Cholesky 分解	301
6.4 带状线性方程组的直接法	303
6.4.1 三对角方程组	303
6.4.2 块三对角方程组	310
6.5 直接法的舍入误差分析	315
6.5.1 矩阵的条件数	315
6.5.2 矩阵条件数的估算	315
6.5.3 舍入误差对解的影响	318
习题 6	319
第 7 章 矩阵特征值问题的数值方法	321
7.1 矩阵的特征值估计和隔离	321
7.2 幂法和反幂法	326
7.2.1 幂法	326
7.2.2 幂法的加速技术	329
7.2.3 反幂法	330
7.3 Jacobi 方法	332
7.3.1 实对称矩阵的旋转正交相似变换	332
7.3.2 Jacobi 方法及其收敛性	335
7.4 QR 方法	338
7.4.1 化一般矩阵为上 Hessenberg 矩阵	339
7.4.2 上 Hessenberg 矩阵的 QR 分解	344
7.4.3 基本 QR 方法	347
7.4.4 带原点位移的 QR 方法	353

7.4.5	双重步位移隐式 QR 方法	354
7.4.6	特征向量的计算方法	361
7.5	Givens-Householder 方法	366
7.5.1	求对称三对角矩阵特征值的二分法	366
7.5.2	二分法的程序实现	371
7.5.3	特征向量的计算	372
7.6	Krylov 子空间方法	374
7.6.1	Rayleigh-Ritz 投影方法	376
7.6.2	Lanczos 方法	379
7.6.3	Arnoldi 方法	396
7.6.4	Jacobi-Davidson 方法	401
	习题 7	405
	参考文献	408

第 1 章 数值线性代数理论基础

本章介绍后面章节中需要用到的一些矩阵代数基础知识: 首先介绍一些基本概念和几种常用的矩阵分解; 其次介绍作为数值误差分析度量工具的向量和矩阵范数及其性质; 接着介绍矩阵的广义逆以及几种特殊的矩阵类型; 最后介绍数值代数中常用的模型问题.

1.1 一些概念和记号

本书引用下列记号. 用 \mathbb{C}^n (\mathbb{R}^n) 表示 n 维复 (实) 向量空间, $\mathbb{C}^{m \times n}$ ($\mathbb{R}^{m \times n}$) 表示 $m \times n$ 阶复 (实) 矩阵空间, $\mathbb{C}_r^{m \times n}$ ($\mathbb{R}_r^{m \times n}$) 表示秩为 r 的 $m \times n$ 阶复 (实) 矩阵集合. 对于任意的矩阵 $A \in \mathbb{C}^{n \times n}$, A^T, A^H, A^{-1} 分别表示矩阵 A 的转置矩阵、共轭转置矩阵和逆矩阵. 对于任意的向量 $x, y \in \mathbb{C}^n$, 用 $(x, y) := y^H x$ 表示 x 和 y 的欧几里得内积.

定义 1.1 设 $A \in \mathbb{C}^{n \times n}$. 若存在数 $\lambda \in \mathbb{C}$ 和非零向量 $x \in \mathbb{C}^n$ 使得 $Ax = \lambda x$, 则称 λ 为 A 的特征值, x 为 A 属于 λ 的特征向量.

由定义 1.1 可知, λ 是 A 的特征值当且仅当 $\det(\lambda I - A) = 0$. 称 $p(\lambda) = \det(\lambda I - A)$ 为 A 的特征多项式.

注意到 $\det(\lambda I - A^T) = \det(\lambda I - A)$, 此即 A 和 A^T 具有相同的特征值. 故存在非零向量 $y \in \mathbb{C}^n$ 满足 $A^T y = \lambda y$, 即 $y^T A = \lambda y^T$, 称 y 为 A 属于 λ 的左特征向量. 相应地, 将 x 称为 A 属于 λ 的右特征向量. 一般来说, 左右特征向量不相等.

如果 A 有 r 个互不相同的特征值 $\lambda_1, \lambda_2, \dots, \lambda_r$, 则 A 的特征多项式可表示为

$$\det(\lambda I - A) = (\lambda - \lambda_1)^{n_1} (\lambda - \lambda_2)^{n_2} \cdots (\lambda - \lambda_r)^{n_r},$$

称 n_i 为 λ_i 的代数重数. 记 $\gamma_i = n - \text{rank}(\lambda_i I - A)$, 称 γ_i 为 λ_i 的几何重数, 它表示属于 λ_i 的线性无关特征向量的个数, 满足 $1 \leq \gamma_i \leq n_i$. 若 $n_i = 1$, 则称 λ_i 为 A 的简单特征值. 若 $\gamma_i = n_i$, 则称 λ_i 为 A 的半简单特征值. 显然, 简单特征值必为半简单特征值. 满足 $p(A) = O$ 的首项系数为 1 且次数最低的多项式 p 称为 A 的最小多项式. 可以证明, A 的最小多项式具有下列形式, 即

$$p(\lambda) = (\lambda - \lambda_1)^{l_1} (\lambda - \lambda_2)^{l_2} \cdots (\lambda - \lambda_r)^{l_r}, \quad 1 \leq l_i \leq n_i, \quad i = 1, 2, \dots, r.$$

定义 1.2 设 $A \in \mathbb{C}^{n \times n}$. 若 $A^H A = I$, 则称 A 为酉矩阵. 实的酉矩阵称为正交矩阵, 即 A 满足 $A^T A = I$.

由定义 1.2 显然有, 酉矩阵 A 的逆矩阵为 $A^{-1} = A^H$, 正交矩阵 A 的逆矩阵为 $A^{-1} = A^T$.

由酉矩阵的定义, 容易证明酉矩阵还具有下列性质.

性质 1.1 设 $U, V \in \mathbb{C}^{n \times n}$ 是酉矩阵, 则 $|\det(U)| = 1$, 且 U^{-1} 和 UV 也是酉矩阵.

性质 1.2 $U \in \mathbb{C}^{n \times n}$ 是酉矩阵的充分必要条件是 U 的列向量是两两正交的单位向量.

定义 1.3 设 $A, B \in \mathbb{C}^{n \times n}$. 若存在非奇异矩阵 P 使得 $A = P^{-1}BP$, 则称 A 与 B 相似. 若 P 为酉矩阵, 则称为酉相似. 若 P 为正交矩阵, 则称为正交相似.

显然, 相似矩阵有相同的特征多项式, 因而有相同的特征值.

定义 1.4 设 $A \in \mathbb{C}^{n \times n}$. 若 $A^H A = A A^H$, 则称 A 为正规矩阵. 若 $A^H = A$, 则称 A 为 Hermite 矩阵. 实的 Hermite 矩阵称为实对称矩阵, 即 A 满足 $A^T = A$. 若 A 满足 $A^H = -A$, 则称 A 为反 Hermite 矩阵. 实的反 Hermite 矩阵称为反对称矩阵, 即 A 满足 $A^T = -A$.

容易验证, 酉矩阵、正交矩阵、Hermite 矩阵、反 Hermite 矩阵、实对称矩阵、反对称矩阵都是正规矩阵.

根据反对称矩阵的定义容易证明反对称矩阵的下列性质.

性质 1.3 设 $A = (a_{ij}) \in \mathbb{R}^{n \times n}$ 是反对称矩阵, 即满足 $A^T = -A$. 则

- (1) $a_{ii} = 0, i = 1, 2, \dots, n$.
- (2) 不存在奇数阶非奇异反对称矩阵.
- (3) A 的特征值只能是 0 或纯虚数.

注 1.1 若 $A \in \mathbb{C}^{n \times n}$ 为 Hermite 矩阵, 则其对角元为 0 或纯虚数, 其特征值也只能是 0 或纯虚数.

定义 1.5 设 $A \in \mathbb{R}^{n \times n}$. 若对任意的非零向量 $x \in \mathbb{R}^n$ 有 $(Ax, x) > 0$ (≥ 0), 则称 A 为 (实) 正定矩阵 (半正定矩阵). 若 A 还是对称的, 则称为 (实) 对称正定矩阵 (对称半正定矩阵).

由定义 1.5 不难推得, 对称矩阵的特征值均为实数. 对称正定矩阵 (对称半正定矩阵) 的特征值均为正数 (非负数). 此外, 若记

$$H = \frac{1}{2}(A + A^T), \quad S = \frac{1}{2}(A - A^T),$$

分别称为 A 的对称部分和反对称部分, 则显然对任意的矩阵 $A \in \mathbb{R}^{n \times n}$ 都可唯一地分裂为

$$A = H + S,$$

称为矩阵 A 的对称-反对称分裂.

性质 1.4 $A \in \mathbb{R}^{n \times n}$ 正定 (或半正定) 的充分必要条件是 its 对称部分 H 对称正定 (或对称半正定).

定义 1.6 设 i_1, i_2, \dots, i_n 是 $1, 2, \dots, n$ 的一个排列, 以 n 阶单位矩阵 I_n 的 n 个列向量

$$e_1 = (1, 0, \dots, 0)^T, e_2 = (0, 1, \dots, 0)^T, \dots, e_n = (0, 0, \dots, 1)^T$$

为列构成的 n 阶矩阵 $P = [e_{i_1}, e_{i_2}, \dots, e_{i_n}]$, 称为置换矩阵或排列矩阵.

例如

$$P = [e_3, e_1, e_2] = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}$$

是一个 3 阶置换矩阵.

性质 1.5 置换矩阵具有如下性质:

- (1) 置换矩阵的转置仍是置换矩阵.
- (2) 置换矩阵是正交矩阵.
- (3) 设 $A \in \mathbb{C}^{n \times n}$, $P = [e_{i_1}, e_{i_2}, \dots, e_{i_n}]$, 则 $P^T A$ 是将 A 按 i_1, i_2, \dots, i_n 行重新排列所得到的矩阵, AP 是将 A 按 i_1, i_2, \dots, i_n 列重新排列得到的矩阵.

定义 1.7 设 $A, B \in \mathbb{C}^{n \times n}$. 如果存在 n 阶非奇异矩阵 C , 使得

- (1) $B = C^T A C$, 则称 A 与 B 为 T-合同.
- (2) $B = C^H A C$, 则称 A 与 B 为 H-合同.

显然, 这两个合同概念具有密切的联系. 当 C 是实矩阵时, T-合同和 H-合同是一致的. 此外, 容易证明, T-合同和 H-合同都是等价关系.

利用合同的定义还可以得到:

- (1) 如果 A 是 Hermite 矩阵, 则 $C^H A C$ 也是 Hermite 矩阵 (即使 C 是奇异矩阵).
- (2) 如果 A 是对称矩阵 (不一定是实矩阵), 则 $C^T A C$ 也是对称矩阵.
- (3) 如果 A 是 Hermite 正定 (半正定) 矩阵, 则 $C^H A C$ 也是 Hermite 正定 (半正定) 矩阵.
- (4) 如果 A 是对称正定 (半正定) 矩阵, 则 $C^T A C$ 也是对称正定 (半正定) 矩阵.

定义 1.8 对于 2×2 分块矩阵

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}, \quad (1.1)$$

当 A_{11} 可逆时, 称 $A_{22} - A_{21}A_{11}^{-1}A_{12}$ 为 A 关于 A_{11} 的 Schur 补, 记为 A/A_{11} . 当 A_{22} 可逆时, 称 $A_{11} - A_{12}A_{22}^{-1}A_{21}$ 为 A 关于 A_{22} 的 Schur 补, 记为 A/A_{22} .

定理 1.1 设 A 具有式 (1.1) 的分块形式, 且 A_{11} 可逆. 则

- (1) $\det(A) = \det(A_{11}) \det(A/A_{11})$.
- (2) $\text{rank}(A) = \text{rank}(A_{11}) + \text{rank}(A/A_{11})$.

证明 注意到

$$\begin{aligned} B &= \begin{bmatrix} I & O \\ -A_{21}A_{11}^{-1} & I \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} I & -A_{11}^{-1}A_{12} \\ O & I \end{bmatrix} \\ &= \begin{bmatrix} A_{11} & O \\ O & A_{22} - A_{21}A_{11}^{-1}A_{12} \end{bmatrix} = \begin{bmatrix} A_{11} & O \\ O & A/A_{11} \end{bmatrix}. \end{aligned}$$

由此立即可得定理的结论. 证毕. □

由定理 1.1 立即可得下面的结论.

推论 1.1 设 A 具有式 (1.1) 的分块形式, 且 A_{22} 可逆, 则

- (1) $\det(A) = \det(A_{22}) \det(A/A_{22})$.
- (2) $\text{rank}(A) = \text{rank}(A_{22}) + \text{rank}(A/A_{22})$.

推论 1.2 设

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

是 Hermite 矩阵, A 关于 A_{11} 的 Schur 补 $A/A_{11} = A_{22} - A_{21}A_{11}^{-1}A_{12}$. 则

- (1) A 正定当且仅当 A_{11} 及 A/A_{11} 均正定.
- (2) 若 A_{11} 正定, 则 A 正定当且仅当 A/A_{11} 正定.

1.2 几种常用的矩阵分解

所谓矩阵分解, 就是将一个矩阵分解为 (从某种意义上讲) 比较简单或对其性质比较熟悉的若干个矩阵的乘积. 下面介绍几种常用的矩阵分解.

1.2.1 矩阵的特征分解

特征分解, 又称谱分解, 是将矩阵分解为由其特征值和特征向量表示的矩阵之积的方法. 注意: 只有对可对角化矩阵才可以实施特征分解.

设 $A \in \mathbb{C}^{n \times n}$ 有 n 个线性无关的特征向量 $q_i (i = 1, 2, \dots, n)$. 则 A 可以被分解为

$$A = Q\Lambda Q^{-1}, \quad (1.2)$$

式中: $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$, $\lambda_i (i = 1, 2, \dots, n)$ 为 A 的特征值; Q 为 n 阶方阵, 且其第 i 列为 A 的特征向量 q_i .

一般来说, 特征向量 $q_i (i = 1, 2, \dots, n)$ 通常被正交单位化 (但这不是必须的). 未被正交单位化的特征向量组 $v_i (i = 1, 2, \dots, n)$ 也可以作为 Q 的列向量.

通常可以通过特征分解来求矩阵的逆. 若矩阵 A 可被特征分解并特征值中不含零, 则矩阵 A 为非奇异矩阵, 且其逆矩阵可以由下式给出:

$$A^{-1} = Q\Lambda^{-1}Q^{-1} = Q\Lambda^{-1}Q^T, \quad (1.3)$$

此处假设 Q 的列向量被正交单位化. 因为 Λ 为对角矩阵, 其逆矩阵容易计算出:

$$\Lambda^{-1} = \text{diag}(\lambda_1^{-1}, \lambda_2^{-1}, \dots, \lambda_n^{-1}).$$

由于任意的 n 阶实对称矩阵 A 都有 n 个线性无关的特征向量, 并且这些特征向量都可以正交单位化而得到一组正交且模为 1 的向量. 故实对称矩阵 A 可被分解成

$$A = Q\Lambda Q^T,$$

式中: Q 为正交矩阵; A 为实对角矩阵.

类似地, 设 $A \in \mathbb{C}^{n \times n}$ 为正规矩阵, 由于正规矩阵是可对角化的 (见定理 1.5), 故其具有一组标准正交特征向量基, 因此可被分解成

$$A = UAU^H,$$

式中: U 为酉矩阵. 进一步地, 若 A 是 Hermite 矩阵, 那么对角矩阵 A 的对角元全为实数. 若 A 是酉矩阵, 则 A 的所有对角元在复平面的单位圆上取得.

下面再给出矩阵特征分解的一个应用.

定理 1.2 设 $A \in \mathbb{C}^{n \times n}$ 为 Hermite 正定 (半正定) 矩阵, 则存在唯一的 Hermite 正定 (半正定) 矩阵 B 使得

$$A = B^2. \quad (1.4)$$

称这样定义的矩阵 B 为矩阵 A 的平方根矩阵, 常记为 $A^{1/2}$.

证明 只证明正定的情形. 存在性. 由于 A 是 Hermite 正定矩阵, 故存在酉矩阵 P 使得

$$A = P \operatorname{diag}(\lambda_1, \lambda_2, \dots, \lambda_n) P^H,$$

式中: $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n > 0$ 为 A 的特征值. 令

$$B = P \operatorname{diag}(\lambda_1^{1/2}, \lambda_2^{1/2}, \dots, \lambda_n^{1/2}) P^H,$$

则显然 B 为 Hermite 正定矩阵, 且 $B^2 = A$.

唯一性. 假设另有 Hermite 正定矩阵 C 满足 $A = C^2$. 则存在酉矩阵 Q 使得

$$C = Q \operatorname{diag}(\mu_1, \mu_2, \dots, \mu_n) Q^H,$$

式中: $\mu_1 \geq \mu_2 \geq \dots \geq \mu_n > 0$ 为 C 的特征值. 由于

$$C^2 = Q \operatorname{diag}(\mu_1^2, \mu_2^2, \dots, \mu_n^2) Q^H = A,$$

故 $\mu_i^2 = \lambda_i$ ($i = 1, 2, \dots, n$), 即 $\mu_i = \lambda_i^{1/2}$. 于是

$$C = Q \operatorname{diag}(\lambda_1^{1/2}, \lambda_2^{1/2}, \dots, \lambda_n^{1/2}) Q^H.$$

再由 $A = B^2 = C^2$, 得

$$P \operatorname{diag}(\lambda_1, \lambda_2, \dots, \lambda_n) P^H = Q \operatorname{diag}(\lambda_1, \lambda_2, \dots, \lambda_n) Q^H,$$

或等价地, 有

$$Q^H P \operatorname{diag}(\lambda_1, \lambda_2, \dots, \lambda_n) = \operatorname{diag}(\lambda_1, \lambda_2, \dots, \lambda_n) Q^H P.$$

记 $Q^H P = (t_{ij})$, 则有

$$t_{ij} \lambda_j = \lambda_i t_{ij}, \quad i, j = 1, 2, \dots, n.$$

因此, 当 $\lambda_i \neq \lambda_j$ 时, 必有 $t_{ij} = 0$. 故无论 λ_i 与 λ_j 是否相等, 都有

$$t_{ij}\lambda_j^{1/2} = t_{ij}\lambda_i^{1/2}, \quad i, j = 1, 2, \dots, n.$$

即

$$Q^H P \operatorname{diag}(\lambda_1^{1/2}, \lambda_2^{1/2}, \dots, \lambda_n^{1/2}) = \operatorname{diag}(\lambda_1^{1/2}, \lambda_2^{1/2}, \dots, \lambda_n^{1/2}) Q^H P,$$

或者等价地, 有

$$P \operatorname{diag}(\lambda_1^{1/2}, \lambda_2^{1/2}, \dots, \lambda_n^{1/2}) P^H = Q \operatorname{diag}(\lambda_1^{1/2}, \lambda_2^{1/2}, \dots, \lambda_n^{1/2}) Q^H,$$

即 $B = C$. 证毕. □

1.2.2 矩阵的 Schur 分解

矩阵的 Schur 分解在理论上十分重要, 它是许多重要定理证明的出发点. 如矩阵论中极为重要的 Hamilton-Cayley (哈密顿-凯莱) 定理, 就可以利用矩阵的 Schur 分解定理进行简洁而优美的证明.

定理 1.3 (Schur 分解定理) 设 $A \in \mathbb{C}^{n \times n}$, 则存在酉矩阵 $P \in \mathbb{C}^{n \times n}$ 使得

$$P^H A P = T,$$

式中: T 为上三角矩阵, 其对角元素是 A 的特征值, 而且可以选取 P 使得 T 的对角元可以任意排列.

证明 对矩阵的阶数 n 用数学归纳法. 当 $n = 1$ 时, 结论显然成立. 假设结论对 $n - 1$ 阶矩阵成立, 下证对 n 阶矩阵成立.

设 A 的特征值为 $\lambda_1, \lambda_2, \dots, \lambda_n$, x_1 为对应于特征值 λ_1 的单位特征向量. 将 x_1 扩充为向量空间 \mathbb{C}^n 的一组标准正交基 x_1, z_2, \dots, z_n . 定义酉矩阵 $P_0 = [x_1, z_2, \dots, z_n]$, 则

$$P_0^H A P_0 = \begin{bmatrix} \lambda_1 & * \\ 0 & A_1 \end{bmatrix},$$

式中: A_1 为 $n - 1$ 阶方阵, 其特征值为 $\lambda_2, \dots, \lambda_n$.

由归纳法假设, 存在 $n - 1$ 阶酉矩阵 P_1 使得

$$P_1^H A_1 P_1 = T_1,$$

式中: T_1 为 $n - 1$ 阶上三角矩阵, 其对角元为 A_1 的特征值 $\lambda_2, \dots, \lambda_n$. 令

$$P = P_0 \begin{bmatrix} 1 & 0 \\ 0 & P_1 \end{bmatrix}, \quad T = \begin{bmatrix} \lambda_1 & * \\ 0 & T_1 \end{bmatrix}.$$

则有

$$P^H A P = \begin{bmatrix} 1 & 0 \\ 0 & P_1^H \end{bmatrix} P_0^H A P_0 \begin{bmatrix} 1 & 0 \\ 0 & P_1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & P_1^H \end{bmatrix} \begin{bmatrix} \lambda_1 & * \\ 0 & A_1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & P_1 \end{bmatrix}$$

$$= \begin{bmatrix} \lambda_1 & * \\ \mathbf{0} & P_1^H A_1 P_1 \end{bmatrix} = \begin{bmatrix} \lambda_1 & * \\ \mathbf{0} & T_1 \end{bmatrix} = T.$$

由归纳法假设, T_1 的对角元可以任意排列, 故 T 的对角元也可以任意排列. 证毕. \square

推论 1.3 设 $A \in \mathbb{R}^{n \times n}$ 的特征值都是实数, 则 A 正交相似于上三角矩阵.

设多项式

$$f(\lambda) = b_0 \lambda^m + b_1 \lambda^{m-1} + \cdots + b_{m-1} \lambda + b_m,$$

对于 n 阶矩阵 A , 定义矩阵多项式

$$f(A) = b_0 A^m + b_1 A^{m-1} + \cdots + b_{m-1} A + b_m I,$$

式中: I 为单位矩阵. 如果 $f(A) = O$, 那么称 A 为 $f(\lambda)$ 的矩阵根.

下面利用矩阵的 Schur 分解来证明著名的 Hamilton-Cayley 定理.

定理 1.4 (Hamilton-Cayley 定理) 设 $A \in \mathbb{C}^{n \times n}$, $p(\lambda) = \det(\lambda I - A)$, 则 $p(A) = O$. 即矩阵 A 的特征多项式是它的零化多项式.

证明 设 n 阶矩阵 A 的特征多项式为

$$p(\lambda) = \det(\lambda I - A) = \lambda^n + a_1 \lambda^{n-1} + \cdots + a_{n-1} \lambda + a_n,$$

记 A 的 n 个特征值为 $\lambda_1, \lambda_2, \cdots, \lambda_n$, 则有

$$p(\lambda) = (\lambda - \lambda_1)(\lambda - \lambda_2) \cdots (\lambda - \lambda_n),$$

由定理 1.3 可知, 存在酉矩阵 U , 使得

$$U^H A U = \begin{bmatrix} \lambda_1 & * & \cdots & * \\ & \lambda_2 & \ddots & \vdots \\ & & \ddots & * \\ & & & \lambda_n \end{bmatrix},$$

于是,

$$p(U^H A U) = (U^H A U - \lambda_1 I)(U^H A U - \lambda_2 I) \cdots (U^H A U - \lambda_n I)$$

$$= \begin{bmatrix} 0 & * & \cdots & * \\ & \lambda_2 - \lambda_1 & \ddots & \vdots \\ & & \ddots & * \\ & & & \lambda_n - \lambda_1 \end{bmatrix} \times \begin{bmatrix} \lambda_1 - \lambda_2 & * & \cdots & * \\ & 0 & \ddots & \vdots \\ & & \ddots & * \\ & & & \lambda_n - \lambda_2 \end{bmatrix}$$

$$\begin{aligned}
& \times \begin{bmatrix} \lambda_1 - \lambda_3 & * & * & \cdots & * \\ & \lambda_2 - \lambda_3 & * & \cdots & * \\ & & 0 & \ddots & \vdots \\ & & & \ddots & * \\ & & & & \lambda_n - \lambda_3 \end{bmatrix} \times \cdots \times \begin{bmatrix} \lambda_1 - \lambda_n & * & * & \cdots & * \\ & \lambda_2 - \lambda_n & * & \cdots & * \\ & & \ddots & \ddots & \vdots \\ & & & \ddots & * \\ & & & & \lambda_{n-1} - \lambda_n \\ & & & & & 0 \end{bmatrix} \\
& = \begin{bmatrix} 0 & 0 & * & \cdots & * \\ 0 & 0 & * & \cdots & * \\ 0 & 0 & * & \cdots & * \\ 0 & 0 & & \ddots & \vdots \\ 0 & 0 & & & * \end{bmatrix} \times \begin{bmatrix} \lambda_1 - \lambda_3 & * & * & \cdots & * \\ & \lambda_2 - \lambda_3 & * & \cdots & * \\ & & 0 & \ddots & \vdots \\ & & & \ddots & * \\ & & & & \lambda_n - \lambda_3 \end{bmatrix} \\
& \times \cdots \times \begin{bmatrix} \lambda_1 - \lambda_n & * & * & \cdots & * \\ & \lambda_2 - \lambda_n & * & \cdots & * \\ & & \ddots & \ddots & \vdots \\ & & & \ddots & * \\ & & & & \lambda_{n-1} - \lambda_n \\ & & & & & 0 \end{bmatrix} = O,
\end{aligned}$$

即 $U^H p(A) U = O$, 也就是 $p(A) = O$. 证毕. \square

在计算数学领域, 人们颇为关心的是, 通过相似变换可以把一个矩阵变换成何种最简单的形状. 利用定理 1.3 和推论 1.3, 可以导出矩阵酉 (正交) 相似于对角矩阵的充要条件.

定理 1.5 矩阵 $A \in \mathbb{C}^{n \times n}$ 酉相似于对角矩阵当且仅当 A 是正规矩阵, 即 $A^H A = A A^H$. 换言之, 正规矩阵一定可对角化.

证明 必要性. 设酉矩阵 U , 使得 $U^H A U = \Lambda$, 其中 Λ 为对角矩阵, 则有

$$\begin{aligned}
A &= U \Lambda U^H, \quad A^H = U \bar{\Lambda} U^H, \\
A^H A &= U \bar{\Lambda} U^H U \Lambda U^H = U \bar{\Lambda} \Lambda U^H = U \Lambda \bar{\Lambda} U^H \\
&= U \Lambda U^H U \bar{\Lambda} U^H = A A^H.
\end{aligned}$$

充分性. 由 Schur 分解定理可知, 存在酉矩阵 U , 使得

$$U^H A U = \begin{bmatrix} t_{11} & t_{12} & \cdots & t_{1n} \\ & t_{22} & \cdots & t_{2n} \\ & & \ddots & \vdots \\ & & & t_{nn} \end{bmatrix} := T.$$

利用条件 $A^H A = A A^H$, 得

$$\begin{aligned} T^H T &= U^H A^H U U^H A U = U^H A^H A U \\ &= U^H A A^H U = U^H A U U^H A^H U = T T^H. \end{aligned}$$

比较上式两端矩阵的对应元素, 可得 $t_{ij} = 0$ ($i < j$), 也就是 $T = \text{diag}(t_{11}, t_{22}, \dots, t_{nn})$, 即 A 酉相似于对角矩阵. 证毕. \square

推论 1.4 设 $A \in \mathbb{R}^{n \times n}$ 的 n 个特征值都是实数, 则 A 正交相似于对角矩阵的充要条件是 $A^T A = A A^T$.

推论 1.5 设 $A \in \mathbb{C}^{n \times n}$ 为正规矩阵, λ 是 A 的特征值, x 是对应的特征向量. 则 $\bar{\lambda}$ 是 A^H 的特征值, $\bar{\lambda}$ 对应的特征向量仍然是 x .

证明 由定理 1.5, 存在 n 阶酉矩阵 U 使得 $U^H A U = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$. 于是有 $U^H A^H U = \text{diag}(\bar{\lambda}_1, \bar{\lambda}_2, \dots, \bar{\lambda}_n)$. 设 $U = [u_1, u_2, \dots, u_n]$, 则上面两式可写为

$$A u_i = \lambda_i u_i, \quad A^H u_i = \bar{\lambda}_i u_i, \quad i = 1, 2, \dots, n.$$

由此可见, 若 λ_i 是 A 的特征值且 u_i 是相应的特征向量时, $\bar{\lambda}_i$ 是 A^H 的特征值且相应的特征向量仍是 u_i . 证毕. \square

推论 1.6 设 $A \in \mathbb{C}^{n \times n}$ 为正规矩阵, λ, μ 是 A 的特征值, x, y 是对应的特征向量. 如果 $\lambda \neq \mu$, 则 x 与 y 正交, 即 $y^H x = 0$.

证明 因为 $Ax = \lambda x$, $Ay = \mu y$. 由推论 1.5 得 $A^H x = \bar{\lambda} x$. 故

$$\bar{\mu} y^H x = (\mu y)^H x = (Ay)^H x = y^H A^H x = \bar{\lambda} y^H x,$$

即 $(\bar{\lambda} - \bar{\mu}) y^H x = 0$. 由题设 $\lambda \neq \mu$, 故必有 $y^H x = 0$. 证毕. \square

推论 1.7 设 $A \in \mathbb{C}^{n \times n}$, 则 A 酉相似于实对角矩阵的充要条件是 A 为 Hermite 矩阵.

证明 必要性. 设存在 n 阶酉矩阵 U 使得

$$U^H A U = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n) = \Lambda,$$

式中: λ_i ($i = 1, 2, \dots, n$) 为实数, 则有

$$A^H = (U \Lambda U^H)^H = U \Lambda^H U^H = U \Lambda U^H = A,$$

即 A 是 Hermite 矩阵.

充分性. 若 A 是 Hermite 矩阵, 则 A 为正规矩阵, 于是存在酉矩阵 U 使得

$$U^H A U = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n) = \Lambda,$$

故

$$A^H = (U^H A U)^H = U^H A^H U = U^H A U = A.$$

从而 $\bar{\lambda}_i = \lambda_i$ ($i = 1, 2, \dots, n$), 即 λ_i 都是实数. 证毕. \square

推论 1.8 设 $A \in \mathbb{C}^{n \times n}$ 是反 Hermite 矩阵, 则存在 n 阶酉矩阵 U 使得

$$U^H A U = \text{diag}(ib_1, ib_2, \dots, ib_n),$$

式中: $b_i (i = 1, 2, \dots, n)$ 为实数.

证明 由于反 Hermite 矩阵是正规矩阵且其特征值为 0 或纯虚数, 由定理 1.5 即得结论. 证毕. \square

下面介绍两个 n 阶实对称矩阵“同时”相似于对角矩阵的问题.

定理 1.6 设 A 和 B 都是 n 阶实对称矩阵, 则存在正交矩阵 P , 使得 $P^T A P$ 和 $P^T B P$ 都是对角矩阵的充要条件是 $AB = BA$.

证明 必要性. 设 n 阶正交矩阵 P 使得

$$P^T A P = \Lambda, \quad P^T B P = \Sigma,$$

式中: Λ 和 Σ 为对角矩阵. 则有

$$\begin{aligned} AB &= P \Lambda P^T P \Sigma P^T = P \Lambda \Sigma P^T \\ &= P \Sigma \Lambda P^T = P \Sigma P^T P \Lambda P^T = BA. \end{aligned}$$

充分性. 设 A 的全体互异特征值为 $\lambda_1, \lambda_2, \dots, \lambda_r$, 相应的重数为 n_1, n_2, \dots, n_r ($n_1 + n_2 + \dots + n_r = n$). 由 A 实对称知, 存在正交矩阵 Q 使得

$$Q^T A Q = \begin{bmatrix} \lambda_1 I_{n_1} & & & \\ & \lambda_2 I_{n_2} & & \\ & & \ddots & \\ & & & \lambda_r I_{n_r} \end{bmatrix},$$

将矩阵 $Q^T B Q$ 划分为如下分块形式

$$Q^T B Q = \begin{bmatrix} B_{11} & B_{12} & \cdots & B_{1r} \\ B_{21} & B_{22} & \cdots & B_{2r} \\ \vdots & \vdots & \ddots & \vdots \\ B_{r1} & B_{r2} & \cdots & B_{rr} \end{bmatrix},$$

式中: B_{ij} 为 $n_i \times n_j$ 阶矩阵 ($i, j = 1, 2, \dots, r$).

由 $AB = BA$ 及 $Q^T Q = I$, 得

$$(Q^T A Q)(Q^T B Q) = (Q^T B Q)(Q^T A Q),$$

即

$$\begin{bmatrix} \lambda_1 B_{11} & \lambda_1 B_{12} & \cdots & \lambda_1 B_{1r} \\ \lambda_2 B_{21} & \lambda_2 B_{22} & \cdots & \lambda_2 B_{2r} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_r B_{r1} & \lambda_r B_{r2} & \cdots & \lambda_r B_{rr} \end{bmatrix} = \begin{bmatrix} \lambda_1 B_{11} & \lambda_2 B_{12} & \cdots & \lambda_r B_{1r} \\ \lambda_1 B_{21} & \lambda_2 B_{22} & \cdots & \lambda_r B_{2r} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_1 B_{r1} & \lambda_2 B_{r2} & \cdots & \lambda_r B_{rr} \end{bmatrix},$$

也就是

$$\lambda_i B_{ij} = \lambda_j B_{ij} \quad (i, j = 1, 2, \cdots, r).$$

当 $i \neq j$ 时, 由 $\lambda_i \neq \lambda_j$ 可得 $B_{ij} = O$ ($i \neq j$; $i, j = 1, 2, \cdots, r$). 于是

$$Q^T B Q = \text{diag}(B_{11}, B_{22}, \cdots, B_{rr}),$$

因为 $Q^T B Q$ 实对称, 所以 B_{ii} ($i = 1, 2, \cdots, r$) 实对称, 那么存在正交矩阵 U_i , 使得 $U_i^T B_{ii} U_i = \Sigma_i$ ($i = 1, 2, \cdots, r$), 其中 Σ_i 为对角矩阵. 令 $U = \text{diag}(U_1, U_2, \cdots, U_r)$, 则 U 为正交矩阵, 从而 $P = QU$ 也为正交矩阵, 且有

$$P^T B P = U^T (Q^T B Q) U = \text{diag}(\Sigma_1, \Sigma_2, \cdots, \Sigma_r),$$

$$P^T A P = U^T (Q^T A Q) U = \text{diag}(\lambda_1 I_{n_1}, \lambda_2 I_{n_2}, \cdots, \lambda_r I_{n_r}).$$

证毕. □

定理 1.6 给出了两个实对称矩阵“同时”正交相似于对角矩阵的充要条件, 下面再给出两个 Hermite 矩阵“同时”合同于对角矩阵的充分条件.

定理 1.7 设 A 和 B 都是 n 阶 Hermite 矩阵, 且 B 为正定矩阵, 则存在可逆矩阵 P , 使得 $P^H A P$ 和 $P^H B P$ 都是对角矩阵.

证明 由 B 正定知, 存在酉矩阵 Q_1 , 使得

$$Q_1^H B Q_1 = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} = \Lambda_1, \quad (\lambda_i > 0),$$

令

$$Q_2 = Q_1 \begin{bmatrix} 1/\sqrt{\lambda_1} & & \\ & \ddots & \\ & & 1/\sqrt{\lambda_n} \end{bmatrix},$$

则 Q_2 可逆, 且有 $Q_2^H B Q_2 = I$. 对于 Hermite 矩阵 $Q_2^H A Q_2$, 存在酉矩阵 Q_3 , 使得

$$Q_3^H (Q_2^H A Q_2) Q_3 = \begin{bmatrix} \mu_1 & & \\ & \ddots & \\ & & \mu_n \end{bmatrix} = \Lambda_2.$$

令 $P = Q_2 Q_3$, 则 P 可逆, 且有

$$P^H A P = A_2, \quad P^H B P = I,$$

式中: A_2 和 I 为对角矩阵. 证毕. □

1.2.3 矩阵的奇异值分解

矩阵的奇异值分解在理论上和实际应用中都十分重要. 特别地, 它已成为信息处理、多元统计分析等工程技术领域中不可缺少的工具. 由 1.2.2 节的 Schur 分解定理可知, 任一方阵用酉相似变换只能约化为上三角矩阵, 不能约化为对角矩阵. 而奇异值分解定理表明: 用两个酉矩阵乘到一个矩阵 A 的两边就可以变为对角矩阵. 奇异值分解常用于奇异的或数值上非常接近奇异的矩阵计算. 它不仅能判断矩阵是否接近奇异, 而且也用于数值求解.

对于任意的 $A \in \mathbb{C}^{m \times n}$, 容易验证:

- (1) $A^H A$ 是 Hermite (半) 正定矩阵.
- (2) 齐次线性方程组 $Ax = 0$ 与 $A^H A x = 0$ 同解.
- (3) $\text{rank}(A^H A) = \text{rank}(A)$.
- (4) $A^H A = O \iff A = O$.

定义 1.9 设 $A \in \mathbb{C}^{m \times n} (r \geq 1)$, 记 Hermite 矩阵 $A^H A$ 的 n 个特征值为 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > \lambda_{r+1} = \dots = \lambda_n = 0$, 则称

$$\sigma_i = \sqrt{\lambda_i}, \quad i = 1, 2, \dots, n$$

为 A 的奇异值.

容易看出, A 的奇异值的个数与 A 的列数相同, A 的正奇异值的个数与 A 的秩相同.

定理 1.8 (奇异值分解定理) 设 $A \in \mathbb{C}^{m \times n}$ 的正奇异值为 $\sigma_1, \sigma_2, \dots, \sigma_r$, 则存在 m 阶酉矩阵 U 和 n 阶酉矩阵 V , 使得

$$A = U \Sigma V^H, \tag{1.5}$$

式中: Σ 为 $m \times n$ 阶对角矩阵, 且

$$\Sigma = \begin{bmatrix} \Sigma_r & O \\ O & O \end{bmatrix} \begin{matrix} r \\ m-r \\ r & n-r \end{matrix}, \tag{1.6}$$

这里 $\Sigma_r = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r)$, $\sigma_i > 0$, $1 \leq i \leq r$, $r \leq \min\{m, n\}$.

证明 对任意的 $A \in \mathbb{C}^{m \times n}$, $A^H A$ 是 n 阶 Hermite 半正定矩阵, 它的 n 个非负特征值记为 $\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2$. 不妨设 $\sigma_1^2 \geq \sigma_2^2 \geq \dots \geq \sigma_r^2 > 0$, $\sigma_{r+1}^2 = \sigma_{r+2}^2 = \dots = \sigma_n^2 = 0$, 这里 $r \leq \min\{m, n\}$. $A^H A$ 的特征分解为

$$V^H (A^H A) V = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2), \tag{1.7}$$

式中: $V = [v_1, v_2, \dots, v_n]$ 为 n 阶酉矩阵.

划分 $V = [V_1, V_2]$, 其中 $V_1 = [v_1, v_2, \dots, v_r]$, $V_2 = [v_{r+1}, v_{r+2}, \dots, v_n]$, 令 $\Sigma_r = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r)$, 则有

$$\begin{aligned} V^H(A^H A)V &= \begin{bmatrix} V_1^H \\ V_2^H \end{bmatrix} (A^H A) \begin{bmatrix} V_1 & V_2 \end{bmatrix} \\ &= \begin{bmatrix} V_1^H A^H A V_1 & V_1^H A^H A V_2 \\ V_2^H A^H A V_1 & V_2^H A^H A V_2 \end{bmatrix}. \end{aligned} \quad (1.8)$$

对比式 (1.7) 和式 (1.8), 得

$$\begin{aligned} V_1^H A^H A V_1 &= \Sigma_r^2, & V_2^H A^H A V_2 &= O, \\ V_1^H A^H A V_2 &= O, & V_2^H A^H A V_1 &= O. \end{aligned} \quad (1.9)$$

由式 (1.9) 的第 1, 2 个等式可知 AV_1 的列互相正交, AV_2 的列为零向量. 取 $m \times r$ 阶矩阵

$$U_1 = AV_1 \Sigma_r^{-1}, \quad (1.10)$$

则 $U_1^H U_1 = I_r$. 再取 $U_2 \in \mathbb{C}^{m \times (n-r)}$, 使得 $U = [U_1, U_2]$ 为酉矩阵. 则由 $AV_2 = O$ 及 $U_2^H AV_1 = U_2^H (U_1 \Sigma_r) = (U_2^H U_1) \Sigma_r = O$, 立即推得

$$U^H AV = \begin{bmatrix} U_1^H AV_1 & U_1^H AV_2 \\ U_2^H AV_1 & U_2^H AV_2 \end{bmatrix} = \begin{bmatrix} \Sigma_r & O \\ O & O \end{bmatrix} = \Sigma.$$

证毕. □

在奇异值分解式 (1.5) 中, U 的第 i 列是 A 的对应于奇异值 σ_i 的左奇异向量, V 的第 i 列是 A 的对应于奇异值 σ_i 的右奇异向量. 从定理的证明过程不难看出, A 的奇异值由 A 唯一确定, 但对应于每个奇异值的奇异向量一般不是唯一的.

注 1.2 若 $A \in \mathbb{R}_r^{m \times n}$, 则定理 1.8 中的 U, V 都可以取为实矩阵, 即 U 和 V 分别为 m 阶和 n 阶的正交矩阵.

两个与矩阵 $A \in \mathbb{C}^{m \times n}$ 有关的重要子空间是其列空间和零空间.

定义 1.10 设矩阵 $A \in \mathbb{C}^{m \times n}$, 分别称

$$\mathcal{R}(A) = \{y : y = Ax, x \in \mathbb{C}^n\} \subset \mathbb{C}^m, \quad \mathcal{N}(A) = \{x : Ax = 0, x \in \mathbb{C}^n\} \subset \mathbb{C}^n$$

为 A 的列空间 (值域) 和零空间 (核).

容易验证: 若将 A 按列划分为 $A = [a_1, a_2, \dots, a_n]$, 则

$$\mathcal{R}(A) = \text{span}\{a_1, a_2, \dots, a_n\}.$$

给定两个向量空间 \mathcal{S}_1 和 \mathcal{S}_2 , 它们的和 \mathcal{S} 是一个子空间, 其每一个向量都是 \mathcal{S}_1 的一个向量与 \mathcal{S}_2 的一个向量之和. 两个子空间的交也是一个子空间. 如果 \mathcal{S}_1 和 \mathcal{S}_2 的交退化为 $\{0\}$, 则 \mathcal{S}_1 和 \mathcal{S}_2 之和称为直和, 表示为 $\mathcal{S} = \mathcal{S}_1 \oplus \mathcal{S}_2$. 当 \mathcal{S} 等于 \mathbb{C}^n 时, 则 \mathbb{C}^n 的每一个向量 x 可被唯一地写成 $x = x_1 + x_2$, 其中 $x_1 \in \mathcal{S}_1, x_2 \in \mathcal{S}_2$.

关于矩阵 $A \in \mathbb{C}^{m \times n}$ 的列空间和零空间有下面的正交分解:

$$\mathcal{R}(A) \oplus \mathcal{N}(A^H) = \mathbb{C}^m, \quad \mathcal{N}(A) \oplus \mathcal{R}(A^H) = \mathbb{C}^n.$$

此外, 当 $AS \subset S$ 时, 称子空间 S 在矩阵 A 之下是不变的. 易见, 矩阵 $A \in \mathbb{C}^{n \times n}$ 的列空间 $\mathcal{R}(A)$ 和零空间 $\mathcal{N}(A)$ 都是 A 的不变子空间. 特别地, 对 A 的任意特征值 λ , 子空间 $\mathcal{N}(A - \lambda I)$ 在 A 下不变. 子空间 $\mathcal{N}(A - \lambda I)$ 称为相应于 λ 的特征子空间, 它包含零向量及 A 的所有相应于 λ 的特征向量.

定理 1.9 在 A 的奇异值分解式 (1.5) 中, 记 U 和 V 的列向量分别为 u_1, u_2, \dots, u_m 和 v_1, v_2, \dots, v_n , 则

(1) A 的秩等于其非零奇异值的个数.

(2) $\|A\|_2 = \max_{1 \leq i \leq n} \sigma_i$. 设 σ_1 为最大奇异值, 其对应的右奇异向量为 v_1 , 则 $\|A\|_2 = \|Av_1\|_2$.

(3) $\mathcal{N}(A) = \text{span}\{v_{r+1}, v_{r+2}, \dots, v_n\}$.

(4) $\mathcal{R}(A) = \text{span}\{u_1, u_2, \dots, u_r\}$.

(5) $A = \sum_{i=1}^r \sigma_i u_i v_i^H$.

证明 沿用定理 1.8 中的记号, 式 (1.5) 可写为

$$A = U \Sigma V^H = [U_1, U_2] \begin{bmatrix} \Sigma_r & O \\ O & O \end{bmatrix} \begin{bmatrix} V_1^H \\ V_2^H \end{bmatrix} = U_1 \Sigma_r V_1^H. \quad (1.11)$$

(1) 由于 U 和 V 都是非奇异矩阵, 故立即有 $\text{rank}(A) = \text{rank}(\Sigma) = r$.

(2) 注意到 U 和 V 都是酉矩阵, 故有

$$\begin{aligned} \|A\|_2 &= \max_{\|x\|_2=1} \|Ax\|_2 = \max_{\|x\|_2=1} \|U \Sigma V^H x\|_2 = \max_{\|y\|_2=1} \|\Sigma y\|_2 \\ &= \max \left\{ [(\sigma_1 y_1)^2 + (\sigma_2 y_2)^2 + \dots + (\sigma_n y_n)^2]^{1/2} : \|y\|_2 = 1 \right\} \\ &= \max_{1 \leq i \leq n} \sigma_i. \end{aligned}$$

设 σ_1 为最大奇异值, 其对应的右奇异向量为 v_1 , 则

$$\|Av_1\|_2 = \|U \Sigma V^H v_1\|_2 = \|\Sigma e_1\|_2 = \sigma_1 = \|A\|_2.$$

(3) 有

$$\mathcal{N}(A) = \{x : Ax = 0\} = \{x : U_1 \Sigma_r V_1^H x = 0\} = \{x : V_1^H x = 0\}$$

$$\begin{aligned}
&= \{x : x = k_{r+1}v_{r+1} + \cdots + k_nv_n, k_i \in \mathbb{C}\} \\
&= \text{span}\{v_{r+1}, v_{r+2}, \cdots, v_n\}.
\end{aligned}$$

(4) 由于

$$\begin{aligned}
\mathcal{R}(A) &= \{y : y = Ax\} = \{y : y = U_1(\Sigma_r V_1^H x)\} \subset \mathcal{R}(U_1), \\
\mathcal{R}(U_1) &= \{y : y = U_1 z\} = \{y : y = (U_1 \Sigma_r V_1^T)(V_1 \Sigma_r^{-1} z)\} \\
&= \{y : y = A(V_1 \Sigma_r^{-1} z)\} \subset \mathcal{R}(A),
\end{aligned}$$

故 $\mathcal{R}(A) = \mathcal{R}(U_1) = \text{span}\{u_1, u_2, \cdots, u_r\}$.

(5) 直接计算式 (1.11), 得

$$\begin{aligned}
A &= [u_1, u_2, \cdots, u_r] \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{bmatrix} \begin{bmatrix} v_1^H \\ \vdots \\ v_r^H \end{bmatrix} \\
&= [\sigma_1 u_1, \sigma_2 u_2, \cdots, \sigma_r u_r] \begin{bmatrix} v_1^H \\ \vdots \\ v_r^H \end{bmatrix} \\
&= \sum_{i=1}^r \sigma_i u_i v_i^H.
\end{aligned}$$

证毕. □

下面的定理表明奇异值可以刻画一个矩阵与更低秩的矩阵之间的距离.

定理 1.10 设 $A \in \mathbb{C}^{m \times n}$, 它的奇异值分解由定理 1.8 给出, 其中 $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_r > 0$ 为其非零奇异值, 则有

$$\min_{\text{rank}(B)=k} \|A - B\|_2 = \|A - A_k\|_2 = \sigma_{k+1}, \quad 0 \leq k \leq r-1, \quad (1.12)$$

式中:

$$A_k = \sum_{i=1}^k \sigma_i u_i v_i^H, \quad k > 0, \quad A_0 = O. \quad (1.13)$$

证明 当 $k=0$ 时, $\|A\|_2 = \sigma_1$, 结论显然成立. 设 $k > 0$, 由 A_k 的定义及 A 的奇异值分解, 有

$$\begin{aligned}
U^H A_k V &= \text{diag}(\sigma_1, \cdots, \sigma_k, 0, \cdots, 0), \\
U^H A V &= \text{diag}(\sigma_1, \cdots, \sigma_k, \sigma_{k+1}, \cdots, \sigma_n).
\end{aligned}$$

于是有

$$U^H (A - A_k) V = \text{diag}(0, \cdots, 0, \sigma_{k+1}, \cdots, \sigma_n).$$

注意到 σ_{k+1} 是 $A - A_k$ 的最大奇异值, 立即有 $\|A - A_k\|_2 = \sigma_{k+1}$. 由于 $\text{rank}(A_k) = k$, 故

$$\min_{\text{rank}(B)=k} \|A - B\|_2 \leq \|A - A_k\|_2 = \sigma_{k+1}.$$

下证任意秩为 k 的 $m \times n$ 阶矩阵 B , 均有 $\|A - B\|_2 \geq \sigma_{k+1}$. 由于 $\text{rank}(B) = k$, 故 B 的零空间 $\mathcal{N}(B) = \{x : Bx = 0\}$ 的维数为 $n - k$, 故存在单位正交向量 x_1, \dots, x_{n-k} 使得 $\mathcal{N}(B) = \text{span}\{x_1, \dots, x_{n-k}\}$. 利用维数定理可知

$$S := \text{span}\{x_1, \dots, x_{n-k}\} \cap \text{span}\{v_1, \dots, v_{k+1}\} \neq \{0\}.$$

在 S 中取一非零单位向量 z , 即 z 满足 $\|z\|_2 = 1$, $Bz = 0$ 及

$$z = \sum_{i=1}^{k+1} \alpha_i v_i, \quad Az = \sum_{i=1}^{k+1} \alpha_i Av_i = \sum_{i=1}^{k+1} \sigma_i \alpha_i u_i.$$

注意到 $\|z\|_2^2 = \sum_{i=1}^{k+1} \alpha_i^2 = 1$, 立即有

$$\|A - B\|_2^2 \geq \|(A - B)z\|_2^2 = \|Az\|_2^2 = \sum_{i=1}^{k+1} \sigma_i^2 \alpha_i^2 \geq \sigma_{k+1}^2.$$

证毕. □

1.2.4 矩阵的极分解和满秩分解

本节讨论矩阵的极分解和满秩分解. 极分解是指将一个复方阵分解为一个酉矩阵和 Hermite 半正定矩阵的乘积. 满秩分解是指将一个非零矩阵分解成列满秩矩阵和行满秩矩阵的乘积, 它是研究广义逆矩阵的重要工具之一. 对于极分解有下面的定理.

定理 1.11 设 $A \in \mathbb{C}^{n \times n}$. 则存在 n 阶酉矩阵 U 和 n 阶 Hermite 半正定矩阵 H , 使得

$$A = HU. \quad (1.14)$$

若 A 是非奇异矩阵, 则上述分解是唯一的, 此时, H 是 Hermite 正定矩阵.

证明 由奇异值分解定理可知, 存在酉矩阵 P 和 Q 使得

$$A = P \begin{bmatrix} \Sigma_r & O \\ O & O \end{bmatrix} Q^H,$$

式中: $\Sigma_r = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r)$, $\sigma_i > 0$ 为 A 的奇异值; $r = \text{rank}(A)$. 将上式改写为

$$A = P \begin{bmatrix} \Sigma_r & O \\ O & O \end{bmatrix} P^H P Q^H,$$

令

$$H = P \begin{bmatrix} \Sigma_r & O \\ O & O \end{bmatrix} P^H, \quad U = PQ^H,$$

容易验证 H 为 Hermite 半正定矩阵, 而 U 为酉矩阵.

下面证明当 A 非奇异时分解是唯一的. 事实上, 若 A 还有极分解 $A = H_1 U_1$, 即

$$A = HU = H_1 U_1,$$

式中: U, U_1 为酉矩阵; H, H_1 为 Hermite 正定矩阵. 注意到

$$AA^H = H^2 = H_1^2,$$

则由定理 1.2 可知 $H_1 = H$. 于是, $U_1 = H_1^{-1}A = H^{-1}A = U$. 证毕. \square

注 1.3 若定理 1.11 中的矩阵 A 是实方阵, 则分解式中的 U 和 H 分别为正交矩阵和对称半正定矩阵.

下面讨论满秩分解.

定义 1.11 设 $A \in \mathbb{C}_r^{m \times n} (r > 0)$, 若存在列满秩矩阵 $F \in \mathbb{C}_r^{m \times r}$ 和行满秩矩阵 $G \in \mathbb{C}_r^{r \times n}$, 使得 $A = FG$, 则称 FG 为 A 的一个满秩分解.

定理 1.12 设 $A \in \mathbb{C}_r^{m \times n} (r > 0)$, 则存在 $F \in \mathbb{C}_r^{m \times r}$ 和 $G \in \mathbb{C}_r^{r \times n}$, 使得 $A = FG$.

证明 利用初等变换, 可将 A 化为阶梯形矩阵

$$B = \begin{bmatrix} G \\ O \end{bmatrix}, \quad (G \in \mathbb{C}_r^{r \times n}),$$

即存在有限个初等矩阵的乘积 $P \in \mathbb{C}^{m \times m}$, 使得 $PA = B$. 由于 P 可逆, 划分

$$P^{-1} = [F, S], \quad (F \in \mathbb{C}_r^{m \times r}, S \in \mathbb{C}^{m \times (m-r)}).$$

于是有

$$A = P^{-1}B = [F, S] \begin{bmatrix} G \\ O \end{bmatrix} = FG.$$

证毕. \square

例 1.1 求下列矩阵的满秩分解:

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

解 由

$$[A, I] = \left[\begin{array}{cccc|ccc} -1 & 0 & 1 & 2 & 1 & 0 & 0 \\ 1 & 2 & -1 & 1 & 0 & 1 & 0 \\ 2 & 2 & -2 & -1 & 0 & 0 & 1 \end{array} \right] \xrightarrow{\text{行}} \left[\begin{array}{cccc|ccc} -1 & 0 & 1 & 2 & 1 & 0 & 0 \\ 0 & 2 & 0 & 3 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & -1 & 1 \end{array} \right],$$

可知

$$P = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & -1 & 1 \end{bmatrix}, \quad P^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ -2 & 1 & 1 \end{bmatrix},$$

于是有

$$F = \begin{bmatrix} 1 & 0 \\ -1 & 1 \\ -2 & 1 \end{bmatrix}, \quad G = \begin{bmatrix} -1 & 0 & 1 & 2 \\ 0 & 2 & 0 & 3 \end{bmatrix}, \quad A = FG.$$

定理 1.12 提供的算法求矩阵的满秩分解时需要计算矩阵 P 及 P^{-1} 的前 r 列, 这在高维情形计算量是巨大的. 下面介绍一种避免求逆矩阵的方法.

定义 1.12 设 $B \in \mathbb{C}_r^{m \times n}$ ($r > 0$) 满足: ① B 的后 $m-r$ 行上的元素都是零; ② B 中有 r 个列 (c_1 列, \dots , c_r 列) 构成 I_m 的前 r 个列. 则称 B 为拟 Hermite 标准形.

定理 1.13 设 $A \in \mathbb{C}_r^{m \times n}$ ($r > 0$) 经过初等行变换化为拟 Hermite 标准形 B , 那么 A 有满秩分解 $A = FG$, 其中 F 是由“ A 的 c_1 列, \dots , c_r 列”构成的矩阵, G 是由“ B 的前 r 行”构成的矩阵.

证明 经过初等行变换将 A 化为拟 Hermite 标准形 B , 等价于存在可逆矩阵 $P \in \mathbb{C}^{m \times m}$, 使得 $PA = B$, 或者 $A = P^{-1}B$.

根据 B 中的列标 c_1, \dots, c_r 构造 n 阶置换矩阵

$$P_1 = [e_{c_1}, \dots, e_{c_r}, e_{c_{r+1}}, \dots, e_{c_n}],$$

并对矩阵 A 和 B 按列分块:

$$A = [a_1, a_2, \dots, a_n], \quad B = [b_1, b_2, \dots, b_n].$$

于是可得

$$\begin{aligned} AP_1 &= [a_{c_1}, \dots, a_{c_r}, a_{c_{r+1}}, \dots, a_{c_n}], \\ BP_1 &= [b_{c_1}, \dots, b_{c_r}, b_{c_{r+1}}, \dots, b_{c_n}] = \begin{bmatrix} I_r & B_{12} \\ O & O \end{bmatrix}. \end{aligned}$$

划分 $P^{-1} = [F, S]$, 其中 $F \in \mathbb{C}_r^{m \times r}$, $S \in \mathbb{C}^{m \times (m-r)}$. 根据定理 1.12 可得满秩分解 $A = FG$, 且 G 是由“ B 的前 r 行”构成的矩阵. 由于

$$AP_1 = P^{-1}(BP_1) = [F, S] \begin{bmatrix} I_r & B_{12} \\ O & O \end{bmatrix} = [F, FB_{12}].$$

所以 F 是由“ AP_1 的前 r 列”构成的矩阵, 也就是由“ A 的 c_1 列, \dots , c_r 列”构成的矩阵. 证毕. \square

例 1.2 求下列矩阵的满秩分解:

$$A = \begin{bmatrix} 2 & 1 & 0 & 2 \\ 0 & 0 & 1 & 2 \\ 2 & 1 & 1 & 4 \end{bmatrix}.$$

解 注意到

$$A \xrightarrow{\text{行}} \begin{bmatrix} 2 & 1 & 0 & 2 \\ 0 & 0 & 1 & 2 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

为拟 Hermite 标准形, 且 $c_1 = 2$, $c_2 = 3$, 故

$$F = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}, \quad G = \begin{bmatrix} 2 & 1 & 0 & 2 \\ 0 & 0 & 1 & 2 \end{bmatrix}, \quad A = FG.$$

1.3 向量和矩阵的范数

在许多实际问题中, 常需对同一线性空间中的向量 (或矩阵) 引入作为它们“大小”的一种度量, 进而比较两个向量 (或矩阵) 的“接近”程度. 引入这种体现其“大小”的量就是范数, 它们在理论与实际应用中都占有重要的地位.

1.3.1 向量内积与向量范数

定义 1.13 设 $x = (x_1, x_2, \dots, x_n)^T \in \mathbb{C}^n$, $y = (y_1, y_2, \dots, y_n)^T \in \mathbb{C}^n$, 称复数

$$(x, y) := y^H x = \sum_{i=1}^n x_i \bar{y}_i$$

为向量 x 和 y 的 (欧几里得) 内积或数量积, 而称

$$\|x\|_2 = (x, x)^{\frac{1}{2}} = \left(\sum_{i=1}^n |x_i|^2 \right)^{\frac{1}{2}} \quad (1.15)$$

为向量 x 的 (欧几里得) 范数或 2-范数.

内积的一个重要性质是

$$(Ax, y) = (x, A^H y), \quad A \in \mathbb{C}^{m \times n}, \quad x \in \mathbb{C}^n, \quad y \in \mathbb{C}^m. \quad (1.16)$$

\mathbb{C}^n 中的欧几里得内积与欧几里得范数还具有下列性质.

性质 1.6 设 $x, y \in \mathbb{C}^n$, 则

- (1) $(x, x) = 0$ 当且仅当 $x = 0$.
- (2) $(\lambda x, y) = \lambda(x, y)$, $(x, \lambda y) = \bar{\lambda}(x, y)$, $\forall \lambda \in \mathbb{C}$.
- (3) $(x, y) = \overline{(y, x)}$.
- (4) $(x_1 + x_2, y) = (x_1, y) + (x_2, y)$.
- (5) Cauchy-Schwarz 不等式: $|(x, y)| \leq \|x\|_2 \cdot \|y\|_2$.
- (6) 三角不等式: $\|x + y\|_2 \leq \|x\|_2 + \|y\|_2$.

由 Schwarz 不等式容易证明

$$\|x\|_2 = \max_{\|y\|_2=1} |y^H x|, \quad x \in \mathbb{C}^n. \quad (1.17)$$

下面给出向量范数的一般定义.

定义 1.14 给定 \mathbb{C}^n 中的某个实值函数 $\mathfrak{N}(x) = \|x\|$. 若对任意的 $x, y \in \mathbb{C}^n$ 有

- (1) $\|x\| \geq 0$, 且 $\|x\| = 0$ 当且仅当 $x = 0$.
- (2) $\|\lambda x\| = |\lambda| \cdot \|x\|$, $\forall \lambda \in \mathbb{C}$.
- (3) $\|x + y\| \leq \|x\| + \|y\|$.

则称 $\|x\|$ 为向量 x 的范数. 定义了范数的线性空间称为赋范线性空间.

常用的向量范数有

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

不难证明, 当 $p \rightarrow \infty$ 时, 对于任意的 $x \in \mathbb{C}^n$, 都有 $\|x\|_p \rightarrow \|x\|_\infty$. 此外, 由于 $\|x\| - \|y\| \leq \|x - y\|$, 故向量范数是 \mathbb{C}^n 中的连续函数.

性质 1.7 (向量范数的连续性) \mathbb{C}^n 中的向量范数 $\|x\|$ 是 \mathbb{C}^n 中的连续函数.

性质 1.8 (向量范数的等价性) 设 $\|x\|$ 和 $\|x\|'$ 是 \mathbb{C}^n 中的任意两种范数, 则成立

$$c_1 \|x\| \leq \|x\|' \leq c_2 \|x\|, \quad x \in \mathbb{C}^n, \quad (1.18)$$

式中: $c_1, c_2 > 0$ 为与向量 x 无关的常数.

定义 1.15 设 $\{x^{(k)} = (x_1^{(k)}, x_2^{(k)}, \dots, x_n^{(k)})^T\}_{k=0}^\infty$ 是 \mathbb{C}^n 中的向量序列, $x = (x_1, x_2, \dots, x_n)^T \in \mathbb{C}^n$. 若

$$\lim_{k \rightarrow \infty} x_i^{(k)} = x_i, \quad i = 1, 2, \dots, n,$$

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

利用范数的等价性可以得到向量序列收敛的充分必要条件.

定理 1.14 $\lim_{k \rightarrow \infty} x^{(k)} = x$ 当且仅当 $\lim_{k \rightarrow \infty} \|x^{(k)} - x\| = 0$, 其中 $\|\cdot\|$ 是 \mathbb{C}^n 中任意的向量范数.

1.3.2 矩阵范数与内积

将 $m \times n$ 阶矩阵 A 看作线性空间 $\mathbb{C}^{m \times n}$ 中的元素, 则完全可以按照定义 1.14 的方式引入矩阵的范数. 其中最常用的是与向量 2-范数相对应的范数

$$\|A\|_F = \left(\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 \right)^{1/2} \quad (1.19)$$

称为矩阵 A 的 Frobenius 范数, 简称 F-范数.

利用奇异值分解定理, 容易证明

$$\|A\|_F = \left(\sum_{i=1}^n \sigma_i^2 \right)^{1/2},$$

式中: $\sigma_1, \sigma_2, \dots, \sigma_n$ 为矩阵 A 奇异值.

类似于向量范数的定义, 可以给出矩阵范数的一般定义.

定义 1.16 给定 $\mathbb{C}^{m \times n}$ 中的某个实值函数 $\mathfrak{N}(A) = \|A\|$. 若对任意的 $A, B \in \mathbb{C}^{m \times n}$ 有

- (1) $\|A\| \geq 0$, 且 $\|A\| = 0$ 当且仅当 $A = O$.
- (2) $\|\lambda A\| = |\lambda| \cdot \|A\|, \forall \lambda \in \mathbb{C}$.
- (3) $\|A + B\| \leq \|A\| + \|B\|$.

则称 $\|A\|$ 为矩阵 A 的范数. 若矩阵范数满足

$$\|Ax\| \leq \|A\| \cdot \|x\|, \quad \forall x \in \mathbb{C}^n, A \in \mathbb{C}^{m \times n},$$

则称矩阵范数 $\|\cdot\|$ 和向量范数 $\|\cdot\|$ 是相容的, 其中 $\|Ax\|$ 和 $\|x\|$ 分别是 \mathbb{C}^m 和 \mathbb{C}^n 中的向量范数.

跟向量范数的等价性一样, 可以证明 $\mathbb{C}^{m \times n}$ 中的任意两个矩阵范数也等价, 即存在与矩阵 A 无关的常数 $c_1, c_2 > 0$ 使得

$$c_1 \|A\| \leq \|A\|' \leq c_2 \|A\|.$$

定义 1.17 设 $A \in \mathbb{C}^{m \times n}$, 给定一种向量范数 $\|\cdot\|$, 定义矩阵范数

$$\|A\| = \max_{\|x\| \neq 0} \frac{\|Ax\|}{\|x\|} = \max_{\|x\|=1} \|Ax\|. \quad (1.20)$$

称 $\|A\|$ 为 $\mathbb{C}^{m \times n}$ 中矩阵 A 的算子范数, 或称由该向量范数诱导出来的矩阵范数.

容易证明, 矩阵的算子范数与相应的向量范数是相容的, 且

$$\|AB\| \leq \|A\| \cdot \|B\|, \quad A \in \mathbb{C}^{m \times n}, B \in \mathbb{C}^{n \times p},$$

其中矩阵范数都是由某个向量范数诱导出来的算子范数.

给定 \mathbb{C}^n 中的向量 p -范数, 可以诱导出相应的矩阵 p -算子范数, 记为 $\|\cdot\|_p$. 下面是三种常用的矩阵算子范数.

定理 1.15 设 $A = (a_{ij}) \in \mathbb{C}^{m \times n}$, 则

$$(1) \|A\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}|, \text{ 称为行和范数或 } \infty\text{-范数.}$$

$$(2) \|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}|, \text{ 称为列和范数或 } 1\text{-范数.}$$

(3) $\|A\|_2 = \sigma_1 = \sqrt{\lambda_{\max}(A^H A)}$, 称为谱范数或 2-范数, 其中 σ_1 是矩阵 A 的最大奇异值, $\lambda_{\max}(A^H A)$ 表示矩阵 $A^H A$ 的最大特征值.

证明 先证明 (1). 记 $\sigma = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}|$. 则对任意满足 $\|x\|_\infty = 1$ 的向量 $x \in \mathbb{C}^n$, 有

$$\begin{aligned} \|Ax\|_\infty &= \max_{1 \leq i \leq m} \left| \sum_{j=1}^n a_{ij} x_j \right| \leq \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}| \cdot |x_j| \\ &\leq \left(\max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}| \right) \|x\|_\infty = \sigma. \end{aligned} \quad (1.21)$$

另外, 不妨设 $\sigma = \sum_{j=1}^n |a_{i_0 j}|$. 取

$$\tilde{x} = (\text{sign}(a_{i_0 1}), \text{sign}(a_{i_0 2}), \dots, \text{sign}(a_{i_0 n}))^T,$$

则 $A \neq O$ 蕴含着 $\|\tilde{x}\|_\infty = 1$. 由 $|a| = \text{sign}(a)a$, 有

$$\|A\tilde{x}\|_\infty = \max_{1 \leq i \leq m} \left| \sum_{j=1}^n a_{ij} \tilde{x}_j \right| \geq \left| \sum_{j=1}^n a_{i_0 j} \tilde{x}_j \right| = \sum_{j=1}^n |a_{i_0 j}| = \sigma. \quad (1.22)$$

结论 (2) 的证明与结论 (1) 类似. 最后给出结论 (3) 的证明. 注意到 $A^H A \in \mathbb{C}^{n \times n}$ 是 Hermite 半正定矩阵, 故其特征值为 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$, 相应的单位正交特征向量为 $\xi_1, \xi_2, \dots, \xi_n \in \mathbb{C}^n$. 因此, 对任意满足 $\|x\|_2 = 1$ 的向量 $x \in \mathbb{C}^n$, 有

$$x = \sum_{i=1}^n c_i \xi_i, \quad \sum_{i=1}^n c_i^2 = 1.$$

于是有

$$\begin{aligned} \mathbf{x}^H \mathbf{A}^H \mathbf{A} \mathbf{x} &= \left(\sum_{i=1}^n c_i \xi_i \right)^H \sum_{j=1}^n c_j \mathbf{A}^H \mathbf{A} \xi_j \\ &= \left(\sum_{i=1}^n c_i \xi_i \right)^H \sum_{j=1}^n c_j \lambda_j \xi_j \\ &= \sum_{i=1}^n c_i^2 \lambda_i \xi_i^H \xi_i = \sum_{i=1}^n c_i^2 \lambda_i \leq \lambda_1. \end{aligned}$$

注意到

$$\|\mathbf{A}\|_2 = \max_{\|\mathbf{x}\|_2=1} \|\mathbf{A}\mathbf{x}\|_2 = \max_{\|\mathbf{x}\|_2=1} \left(\mathbf{x}^H \mathbf{A}^H \mathbf{A} \mathbf{x} \right)^{1/2},$$

故 $\|\mathbf{A}\|_2 \leq \sqrt{\lambda_1}$. 取 $\mathbf{x} = \xi_1$, 则有

$$\mathbf{x}^H \mathbf{A}^H \mathbf{A} \mathbf{x} = \xi_1^H \mathbf{A}^H \mathbf{A} \xi_1 = \lambda_1 \xi_1^H \xi_1 = \lambda_1,$$

故

$$\begin{aligned} \|\mathbf{A}\|_2 &= \max_{\|\mathbf{x}\|_2=1} \|\mathbf{A}\mathbf{x}\|_2 = \max_{\|\mathbf{x}\|_2=1} \left(\mathbf{x}^H \mathbf{A}^H \mathbf{A} \mathbf{x} \right)^{1/2} \\ &\geq (\xi_1^H \mathbf{A}^H \mathbf{A} \xi_1)^{1/2} = (\lambda_1 \xi_1^H \xi_1)^{1/2} = \sqrt{\lambda_1}. \end{aligned}$$

因此 $\|\mathbf{A}\|_2 = \sqrt{\lambda_1} = \sqrt{\lambda_{\max}(\mathbf{A}^H \mathbf{A})}$. 证毕. \square

与矩阵的行和范数及列和范数相比, 矩阵的谱范数显得不易计算, 但它有良好的分析性质.

定理 1.16 设 $\mathbf{A} \in \mathbb{C}^{m \times n}$, 则

- (1) $\|\mathbf{A}\|_2 = \max_{\|\mathbf{x}\|_2=\|\mathbf{y}\|_2=1} |\mathbf{y}^H \mathbf{A} \mathbf{x}|$.
- (2) $\|\mathbf{A}^H\|_2 = \|\mathbf{A}^T\|_2 = \|\mathbf{A}\|_2$.
- (3) $\|\mathbf{A}^H \mathbf{A}\|_2 = \|\mathbf{A}\|_2^2$.
- (4) $\|\mathbf{A}\|_2^2 \leq \|\mathbf{A}\|_1 \cdot \|\mathbf{A}\|_\infty$.
- (5) $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F$.
- (6) 设 $\mathbf{U} \in \mathbb{C}^{m \times m}$ 和 $\mathbf{V} \in \mathbb{C}^{n \times n}$ 满足 $\mathbf{U}^H \mathbf{U} = \mathbf{I}_m$ 和 $\mathbf{V}^H \mathbf{V} = \mathbf{I}_n$, 则 $\|\mathbf{U} \mathbf{A} \mathbf{V}\|_2 = \|\mathbf{A}\|_2$.

证明 由算子范数的定义, 存在 $\mathbf{x}_0 \in \mathbb{C}^n$ 满足 $\|\mathbf{x}_0\|_2 = 1$ 和 $\|\mathbf{A}\|_2 = \|\mathbf{A}\mathbf{x}_0\|_2$. 由式 (1.17) 知,

$$\|\mathbf{A}\mathbf{x}_0\|_2 = \max_{\|\mathbf{y}\|_2=1} |\mathbf{y}^H \mathbf{A} \mathbf{x}_0|,$$

这意味着成立

$$\|A\|_2 = \|Ax_0\|_2 \leq \max_{\|x\|_2=1, \|y\|_2=1} |y^H Ax|.$$

另外, 由 Schwarz 不等式及范数的相容性, 有

$$|y^H Ax| \leq \|y\|_2 \cdot \|Ax\|_2 \leq \|y\|_2 \cdot \|A\|_2 \cdot \|x\|_2,$$

故

$$\max_{\|x\|_2=1, \|y\|_2=1} |y^H Ax| \leq \|A\|_2,$$

这就证明了结论 (1).

由结论 (1) 有

$$\begin{aligned} \|A^T\|_2 &= \max_{\|x\|_2=\|y\|_2=1} |y^H A^T x| = \max_{\|x\|_2=\|y\|_2=1} |\bar{x}^H A y| \\ &= \max_{\|\bar{x}\|_2=\|y\|_2=1} |\bar{x}^H A y| = \|A\|_2, \end{aligned}$$

式中: \bar{x}, \bar{y} 为 x, y 的共轭向量, 这里用到了一个标量与其转置相等.

同理, 有

$$\|A^H\|_2 = \max_{\|x\|_2=\|y\|_2=1} |y^H A^H x| = \max_{\|x\|_2=\|y\|_2=1} |x^H A y| = \|A\|_2,$$

这里用到了一个复数的模与其共轭转置的模相等. 至此证明了结论 (2).

注意到 $A^H A$ 是 Hermite 矩阵, 故由谱范数的定义, 有

$$\begin{aligned} \|A^H A\|_2 &= \sqrt{\lambda_{\max}((A^H A)^H (A^H A))} = \sqrt{\lambda_{\max}((A^H A)^2)} \\ &= \lambda_{\max}(A^H A) = \|A\|_2^2, \end{aligned}$$

结论 (3) 得证.

由于 $\|A\|_2^2$ 是 $A^H A$ 的最大特征值, 设其相应的特征向量为 $x \neq 0$, 即 $A^H A x = \|A\|_2^2 x$, 两边取 1-范数, 得

$$\begin{aligned} \|A\|_2^2 \|x\|_1 &= \|\|A\|_2^2 x\|_1 = \|A^H A x\|_1 \\ &\leq \|A^H\|_1 \cdot \|A\|_1 \cdot \|x\|_1 = \|A\|_\infty \cdot \|A\|_1 \cdot \|x\|_1, \end{aligned}$$

两边除以 $\|x\|_1$ 即得结论 (4).

由于 $A^H A$ 是 Hermite 半正定矩阵, 故其特征值 $\lambda_s(A^H A)$, $s = 1, 2, \dots, n$ 均非负, 且满足

$$\sum_{s=1}^n \lambda_s(A^H A) = \text{tr}(A^H A) = \sum_{s=1}^n (A^H A)_{ss},$$

从而有

$$\|A\|_2^2 = \lambda_{\max}(A^H A) \leq \sum_{s=1}^n \lambda_s(A^H A)$$

$$= \sum_{s=1}^n (\mathbf{A}^H \mathbf{A})_{ss} = \sum_{s=1}^n \left(\sum_{i=1}^m |a_{is}|^2 \right) = \|\mathbf{A}\|_F^2,$$

两边开平方即得结论 (5).

最后证明结论 (6). 注意到对任意的 n 阶酉矩阵 \mathbf{U} 和 n 维列向量 \mathbf{x} , 有

$$\|\mathbf{U}\mathbf{x}\|_2 = \sqrt{(\mathbf{U}\mathbf{x})^H(\mathbf{U}\mathbf{x})} = \sqrt{\mathbf{x}^H(\mathbf{U}^H\mathbf{U})\mathbf{x}} = \sqrt{\mathbf{x}^H\mathbf{x}} = \|\mathbf{x}\|_2,$$

即酉矩阵保持列向量的 2-范数不变. 由此, 有

$$\begin{aligned} \|\mathbf{U}\mathbf{A}\mathbf{V}\|_2 &= \max_{\|\mathbf{x}\|_2=1} \|\mathbf{U}\mathbf{A}\mathbf{V}\mathbf{x}\|_2 = \max_{\|\mathbf{y}\|_2=1} \|\mathbf{U}\mathbf{A}\mathbf{y}\|_2 \\ &= \max_{\|\mathbf{y}\|_2=1} \|\mathbf{A}\mathbf{y}\|_2 = \|\mathbf{A}\|_2. \end{aligned} \quad (1.23)$$

证毕. □

定义 1.18 设 $\mathbf{A} \in \mathbb{C}^{n \times n}$, 其特征值为 $\lambda_1, \lambda_2, \dots, \lambda_n$, 则称

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

为矩阵 \mathbf{A} 的谱半径.

由上述定义, $\|\mathbf{A}\|_2$ 可定义为

$$\|\mathbf{A}\|_2 = \sqrt{\rho(\mathbf{A}^H \mathbf{A})}.$$

特别地, 当 \mathbf{A} 为 Hermite 矩阵时, 有

$$\|\mathbf{A}\|_2 = \rho(\mathbf{A}).$$

对于一般情况, 有如下定理.

定理 1.17 设 $\mathbf{A} \in \mathbb{C}^{n \times n}$, 则

(1) \mathbf{A} 的谱半径不超过 \mathbf{A} 的任何范数, 即

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

(2) 对任意正数 ε , 存在 $\mathbb{C}^{n \times n}$ 上的某种算子范数 $\|\cdot\|_\varepsilon$, 使得

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

证明 (1) 设 λ 是 \mathbf{A} 的一个特征值, 且 \mathbf{A} 的谱半径 $\rho(\mathbf{A}) = |\lambda|$, \mathbf{x}_0 是对应于 λ 的特征向量, 即 $\mathbf{A}\mathbf{x}_0 = \lambda\mathbf{x}_0$, 所以对任何一种向量范数, 有 $\|\mathbf{A}\mathbf{x}_0\| = |\lambda| \cdot \|\mathbf{x}_0\|$, 故有

$$\rho(\mathbf{A}) = |\lambda| = \frac{\|\mathbf{A}\mathbf{x}_0\|}{\|\mathbf{x}_0\|} \leq \max_{\|\mathbf{x}\| \neq 0} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} = \|\mathbf{A}\|.$$

(2) 由 Jordan 分解定理知, 存在非奇异矩阵 P 使得

$$P^{-1}AP = \begin{bmatrix} \lambda_1 & \delta_1 & & & \\ & \lambda_2 & \delta_2 & & \\ & & \ddots & \ddots & \\ & & & \lambda_{n-1} & \delta_{n-1} \\ & & & & \lambda_n \end{bmatrix},$$

式中: δ_i 等于 1 或 0. 对于任意的正数 ε , 令

$$D_\varepsilon = \begin{bmatrix} 1 & & & & \\ & \varepsilon & & & \\ & & \ddots & & \\ & & & \varepsilon^{n-1} & \\ & & & & \end{bmatrix},$$

则

$$D_\varepsilon^{-1}P^{-1}APD_\varepsilon = \begin{bmatrix} \lambda_1 & \varepsilon\delta_1 & & & \\ & \lambda_2 & \varepsilon\delta_2 & & \\ & & \ddots & \ddots & \\ & & & \lambda_{n-1} & \varepsilon\delta_{n-1} \\ & & & & \lambda_n \end{bmatrix}.$$

现定义一种新的矩阵范数

$$\|A\|_\varepsilon = \|D_\varepsilon^{-1}P^{-1}APD_\varepsilon\|_\infty, \quad A \in \mathbb{C}^{n \times n},$$

则容易验证这样定义的矩阵范数 $\|\cdot\|_\varepsilon$ 是由如下定义的向量范数诱导出来的算子范数:

$$\|x\|_{PD_\varepsilon} = \|(PD_\varepsilon)^{-1}x\|_\infty, \quad x \in \mathbb{C}^n.$$

从而有

$$\|A\|_\varepsilon = \|D_\varepsilon^{-1}P^{-1}APD_\varepsilon\|_\infty = \max_{1 \leq i \leq n} (|\lambda_i| + \varepsilon\delta_i) \leq \rho(A) + \varepsilon,$$

这里假设 $\delta_n = 0$. 证毕. □

类似于向量序列的极限定义, 有以下定义.

定义 1.19 设矩阵序列 $\{A^{(k)} = (a_{ij}^{(k)}) \in \mathbb{C}^{m \times n}\}_{k=0}^\infty$, $A = (a_{ij}) \in \mathbb{C}^{m \times n}$. 若

$$\lim_{k \rightarrow \infty} a_{ij}^{(k)} = a_{ij}, \quad i = 1, 2, \dots, m; \quad j = 1, 2, \dots, n,$$

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

利用矩阵范数的等价性, 可得以下定理.

定理 1.18 $\lim_{k \rightarrow \infty} A^{(k)} = A$ 当且仅当 $\lim_{k \rightarrow \infty} \|A^{(k)} - A\| = 0$, 其中 $\|\cdot\|$ 是 $\mathbb{C}^{m \times n}$ 中任意的矩阵范数.

关于矩阵范数, 还有下述定理.

定理 1.19 若 $\|A\| < 1$, 则矩阵 $I - A$ 非奇异, 且满足

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

证明 用反证法. 设 $\det(I - A) = 0$, 则方程 $(I - A)x = 0$ 有非零解, 即存在 $x_0 \in \mathbb{C}^n$, $x_0 \neq 0$, 使得 $(I - A)x_0 = 0$. 故

$$\|A\| = \max_{\|x\| \neq 0} \frac{\|Ax\|}{\|x\|} \geq \frac{\|Ax_0\|}{\|x_0\|} = 1.$$

这与 $\|A\| < 1$ 矛盾.

进一步, 由于 $(I - A)(I - A)^{-1} = I$, 则 $(I - A)^{-1} = I + A(I - A)^{-1}$, 从而

$$\|(I - A)^{-1}\| \leq \|I\| + \|A\| \cdot \|(I - A)^{-1}\|.$$

将上式整理即得要证的结论. 证毕. □

定理 1.20 设 $A \in \mathbb{C}^{n \times n}$, 则

$$(1) \quad \lim_{k \rightarrow \infty} A^k = O \iff \rho(A) < 1.$$

$$(2) \quad \lim_{k \rightarrow \infty} \|A^k\|^{\frac{1}{k}} = \rho(A).$$

证明 (1) 充分性. 设 $\rho(A) < 1$, 则存在 $\varepsilon > 0$ 满足 $\rho(A) + 2\varepsilon < 1$. 于是由定理 1.17 知存在算子范数 $\|\cdot\|$ 使得 $\|A\| \leq \rho(A) + \varepsilon$. 因此

$$\|A^k\| \leq \|A\|^k \leq [\rho(A) + \varepsilon]^k < (1 - \varepsilon)^k \rightarrow 0.$$

必要性. 设 $\lim_{k \rightarrow \infty} A^k = O$, 并假定 $\lambda \in \lambda(A)$ 满足 $\rho(A) = |\lambda|$. 由于 $\lambda \in \lambda(A)$ 蕴含 $\lambda^k \in \lambda(A^k)$, 故由定理 1.17 知对一切的 $k = 1, 2, \dots$ 成立

$$[\rho(A)]^k = |\lambda|^k = |\lambda^k| \leq \rho(A^k) \leq \|A^k\|,$$

故必有 $\rho(A) < 1$.

(2) 由结论 (1) 的证明可知, 对任意的 $\varepsilon > 0$, 存在算子范数 $\|\cdot\|$, 满足

$$[\rho(A)]^k \leq \|A^k\| \leq [\rho(A) + \varepsilon]^k,$$

即

$$\rho(A) \leq \|A^k\|^{\frac{1}{k}} \leq \rho(A) + \varepsilon,$$

由 $\varepsilon > 0$ 的任意性, 上式令 $k \rightarrow \infty$ 得

$$\lim_{k \rightarrow \infty} \|A^k\|^{\frac{1}{k}} = \rho(A).$$

再根据有限维空间范数的等价性, 对 $\mathbb{C}^{n \times n}$ 中的任意范数, 上述极限成立. 证毕. □

利用定理 1.20 不难证明下面的重要结论.

定理 1.21 设 $A \in \mathbb{C}^{n \times n}$, 则

(1) 矩阵级数 $\sum_{k=0}^{\infty} A^k$ 收敛的充分必要条件是 $\rho(A) < 1$.

(2) 若矩阵级数 $\sum_{k=0}^{\infty} A^k$ 收敛, 则

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

且对于 $\mathbb{C}^{n \times n}$ 中的任意矩阵范数都有

$$\left\| (I - A)^{-1} - \sum_{k=0}^m A^k \right\| \leq \frac{\|A\|^{m+1}}{1 - \|A\|}.$$

1.4 矩阵的广义逆

矩阵的广义逆是研究最小二乘问题及矩阵方程的一个重要工具. 首先给出矩阵广义逆的定义.

定义 1.20 设 $A \in \mathbb{C}^{m \times n}$, 若有 $X \in \mathbb{C}^{n \times m}$ 满足 Penrose (彭罗斯) 方程, 即

(1) $AXA = A.$

(2) $XAX = X.$

(3) $(AX)^H = AX.$

(4) $(XA)^H = XA.$

则称 X 为 A 的 Moore-Penrose 逆, 记为 A^\dagger .

注 1.4 若 $A \in \mathbb{R}^{m \times n}$, 则 $X = A^\dagger \in \mathbb{R}^{n \times m}$ 定义为满足下列四个方程的 X :

(1) $AXA = A.$

(2) $XAX = X.$

(3) $(AX)^T = AX.$

(4) $(XA)^T = XA.$

由定义 1.20 容易得到:

(1) $A \in \mathbb{C}^{n \times n}$ 可逆时, $X = A^{-1}$ 满足 Penrose 方程, 故 $A^\dagger = A^{-1}$.

(2) $A = O_{m \times n}$ 时, $X = O_{n \times m}$ 满足 Penrose 方程, 故 $O_{m \times n}^\dagger = O_{n \times m}$.

(3) $A = x \in \mathbb{C}^{n \times 1}$ 时, $X = \frac{1}{\|x\|_2^2} x^H$ 满足 Penrose 方程, 故 $x^\dagger = \frac{1}{\|x\|_2^2} x^H$. 例如:

$$\left[\begin{array}{c} 1 \\ 1 + 2i \end{array} \right]^\dagger = \frac{1}{6} (1, 1 - 2i).$$

定理 1.22 设 $F \in \mathbb{C}_r^{m \times r}$, $G \in \mathbb{C}_r^{r \times n}$ ($r \geq 1$), 则有

$$(1) F^\dagger = (F^H F)^{-1} F^H, \quad F^\dagger F = I_r.$$

$$(2) G^\dagger = G^H (G G^H)^{-1}, \quad G G^\dagger = I_r.$$

证明 先证结论 (1). 对于 F , 记 $X = (F^H F)^{-1} F^H$, 则有

$$F X F = F (F^H F)^{-1} F^H F = F, \quad X F X = (F^H F)^{-1} F^H F X = X,$$

$$(F X)^H = X^H F^H = F (F^H F)^{-1} F^H = F X, \quad (X F)^H = I_r^H = I_r = X F.$$

所以 $F^\dagger = (F^H F)^{-1} F^H$, 并且 $F^\dagger F = I_r$.

同理, 可证结论 (2). 证毕. \square

定理 1.23 设 $A \in \mathbb{C}^{m \times n}$, 则 A^\dagger 存在且唯一.

证明 存在性. 当 $A = O$ 时, $A^\dagger = O$; 当 $A \neq O$ 时, $\text{rank}(A) = r \geq 1$, 从而 A 有满秩分解 $A = FG$ ($F \in \mathbb{C}_r^{m \times r}$, $G \in \mathbb{C}_r^{r \times n}$). 记 $X = G^\dagger F^\dagger$, 则有

$$A X A = F G G^\dagger F^\dagger F G = F (G G^\dagger) (F^\dagger F) G = F G = A,$$

$$X A X = G^\dagger F^\dagger F G G^\dagger F^\dagger = G^\dagger (F^\dagger F) (G G^\dagger) F^\dagger = G^\dagger F^\dagger = X,$$

$$(A X)^H = (F G G^\dagger F^\dagger)^H = (F F^\dagger)^H = F F^\dagger = F G G^\dagger F^\dagger = A X,$$

$$(X A)^H = (G^\dagger F^\dagger F G)^H = (G^\dagger G)^H = G^\dagger G = G^\dagger F^\dagger F G = X A,$$

由定义可得

$$A^\dagger = G^\dagger F^\dagger = G^H (F^H A G^H)^{-1} F^H. \quad (1.24)$$

唯一性. 对于 A , 如果 X 和 Y 都满足 Penrose 方程, 则有

$$\begin{aligned} X &= X A X = X \cdot A Y A \cdot X = X \cdot (A Y) (A X) \\ &= X \cdot (A Y)^H (A X)^H = X \cdot (A X A Y)^H = X \cdot (A Y)^H \\ &= X A Y = X \cdot A Y A \cdot Y = (X A) (Y A) \cdot Y \\ &= (X A)^H (Y A)^H \cdot Y = (Y A X A)^H \cdot Y \\ &= (Y A)^H \cdot Y = Y A Y = Y, \end{aligned}$$

所以 A^\dagger 唯一. 证毕. \square

定理 1.23 的存在性证明实际上给出了使用矩阵满秩分解的方法来计算非零矩阵广义逆. 下面给出一个例子.

例 1.3 已知

$$A = \begin{bmatrix} 1 & 2 & 1 & 2 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 2 & 1 & 2 & 1 \end{bmatrix},$$

求 A 的满秩分解和 A^\dagger .

$$\text{解 } A \xrightarrow{\text{初等行变换}} \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} : c_1 = 1, c_2 = 2,$$

$$A = FG : F = \begin{bmatrix} 1 & 2 \\ 0 & 1 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}, \quad G = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix}.$$

因此, 有

$$F^\dagger = (F^T F)^{-1} F^T = \begin{bmatrix} 6 & 4 \\ 4 & 6 \end{bmatrix}^{-1} F^T = \frac{1}{10} \begin{bmatrix} -1 & -2 & 3 & 4 \\ 4 & 3 & -2 & -1 \end{bmatrix},$$

$$G^\dagger = G^T (G G^T)^{-1} = G^T \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}^{-1} = \frac{1}{2} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix},$$

$$A^\dagger = G^\dagger F^\dagger = \frac{1}{20} \begin{bmatrix} -1 & -2 & 3 & 4 \\ 4 & 3 & -2 & -1 \\ -1 & -2 & 3 & 4 \\ 4 & 3 & -2 & -1 \end{bmatrix}.$$

下面的定理说明使用矩阵奇异值分解的方法来计算广义逆更为便捷.

定理 1.24 设 $A \in \mathbb{C}_r^{m \times n} (r \geq 1)$ 的奇异值分解为

$$A = U \begin{bmatrix} \Sigma_r & O \\ O & O \end{bmatrix} V^H,$$

式中: U 和 V 及 Σ_r 的意义同式 (1.5), 则

$$A^\dagger = V \begin{bmatrix} \Sigma_r^{-1} & O \\ O & O \end{bmatrix} U^H. \quad (1.25)$$

证明 根据 Moore-Penrose 广义逆的定义, 直接验证即可证明.

注 1.5 需要指出, 可逆矩阵的逆 A^{-1} 具有的性质, 对于一般矩阵的广义逆 A^\dagger 不一定具有, 例如:

$$(1) (AB)^\dagger \neq B^\dagger A^\dagger : \text{取 } A = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, B = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}, \text{ 则}$$

$$AB = \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix}, \quad (AB)^\dagger = \frac{1}{2} \begin{bmatrix} 1 & 0 \\ 1 & 0 \end{bmatrix},$$

$$A^\dagger = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad B^\dagger = B^{-1} = \begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix}, \quad B^\dagger A^\dagger = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}.$$

(2) $A^\dagger A \neq AA^\dagger$: A 不是方阵时, 这是明显的. A 是方阵但不可逆时, 取

$$A = \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix}, \quad A^\dagger = \frac{1}{2} \begin{bmatrix} 1 & 0 \\ 1 & 0 \end{bmatrix},$$

则

$$A^\dagger A = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad AA^\dagger = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}.$$

1.5 几种特殊的矩阵类型

1. 不可约矩阵和对角占优矩阵

定义 1.21 设 A 是一个 n ($n \geq 2$) 阶矩阵, 如果存在 n 阶置换矩阵 P , 使得

$$P^T A P = \begin{bmatrix} A_{11} & A_{12} \\ O & A_{22} \end{bmatrix},$$

式中: A_{11} 和 A_{22} 分别为 r 阶和 $n-r$ 阶方阵 ($1 \leq r \leq n-1$). 则称 A 为可约矩阵. 如果不存在这样的置换矩阵, 则称 A 为不可约矩阵.

可约矩阵的一个等价定义如下.

定义 1.21' 设 $A = (a_{ij}) \in \mathbb{C}^{n \times n}$ ($n \geq 2$). 又设 $\mathcal{I} = \{1, 2, \dots, n\}$. 若存在 \mathcal{I} 的一个非空真子集 \mathcal{K} , 使得 $a_{ij} = 0$ ($i \notin \mathcal{K}, j \in \mathcal{K}$), 则称 A 为可约矩阵. 否则, 称 A 为不可约矩阵.

用定义 1.21' 判定一个矩阵是否可约更为方便.

例 1.4 设 $A = (a_{ij})$ 为三对角矩阵, 即满足 $a_{ij} = 0$ ($|i-j| > 1$). 证明: 若

$$a_{i+1,i} \neq 0, \quad a_{i,i+1} \neq 0, \quad i = 1, 2, \dots, n-1,$$

则 A 是不可约矩阵.

证明 用反证法. 若 A 可约, 则由定义 1.21' 知, 存在 \mathcal{I} 的非空真子集 \mathcal{K} , 使得 $a_{ij} = 0$ ($i \notin \mathcal{K}, j \in \mathcal{K}$). 由于 \mathcal{K} 非空, 则 \mathcal{K} 至少有一个数, 设为 r . 由 $a_{r,r+1} \neq 0$ 知, $r+1 \in \mathcal{K}$. 类似可得 $r+2 \in \mathcal{K}, \dots, n \in \mathcal{K}$. 另外, 又由 $a_{r-1,r} \neq 0$ 知, $r-1 \in \mathcal{K}$. 类似可得 $r-2 \in \mathcal{K}, \dots, 1 \in \mathcal{K}$. 即有 $\mathcal{K} = \mathcal{I}$, 这与 \mathcal{K} 是 \mathcal{I} 的真子集矛盾. 故 A 为不可约矩阵.

定义 1.22 设 A 是一个 n 阶矩阵.

(1) 若 A 满足

$$|a_{ii}| \geq \sum_{j \neq i} |a_{ij}|, \quad i = 1, 2, \dots, n,$$

且其中至少有个严格不等式成立, 则称 A 为 (行) 弱对角占优.

(2) 若 A 满足

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}|, \quad i = 1, 2, \dots, n,$$

则称 A 为 (行) 严格对角占优.

类似可定义 (列) 弱对角占优和 (列) 严格对角占优.

以后往往不严格区分按行 (严格) 对角占优或按列 (严格) 对角占优而模糊地统称为 (严格) 对角占优. 弱对角占优且不可约的矩阵简称为不可约对角占优.

引理 1.1 若 n 阶矩阵 A 满足下列条件之一:

(1) 严格对角占优.

(2) 不可约对角占优.

则 A 非奇异, 即其行列式 $\det(A) \neq 0$.

证明 首先考虑 A 为严格对角占优的情形. 用反证法. 假设 A 是奇异的, 则存在非零向量 x 使得 $Ax = 0$. 不妨设 $|x_i| = \|x\|_\infty = 1$. 则有

$$|a_{ii}| = |a_{ii}x_i| = \left| \sum_{j \neq i} a_{ij}x_j \right| \leq \sum_{j \neq i} |a_{ij}|,$$

这与 A 严格对角占优矛盾.

再考虑 A 为不可约对角占优时结论成立. 仍用反证法. 设非零向量 x 满足 $Ax = 0$, 并定义

$$\mathcal{K} = \{k : |x_k| < 1\}.$$

易知 \mathcal{K} 非空. 否则 x 的所有分量的模均为 1, 则对于所有的 $i \in \mathcal{K}$, 均有

$$|a_{ii}| \leq \sum_{j \neq i} |a_{ij}|,$$

这与 A 对角占优矛盾. 此外, 因 A 不可约, 必存在 i, k 使得

$$a_{ik} \neq 0, \quad i \notin \mathcal{K}, \quad k \in \mathcal{K}.$$

于是 $|a_{ik}x_k| < |a_{ik}|$, 且

$$|a_{ii}| \leq \sum_{j \notin \mathcal{K}, j \neq i} |a_{ij}| \cdot |x_j| + \sum_{j \in \mathcal{K}} |a_{ij}| \cdot |x_j| < \sum_{j \notin \mathcal{K}, j \neq i} |a_{ij}| + \sum_{j \in \mathcal{K}} |a_{ij}| = \sum_{j \neq i} |a_{ij}|,$$

这又与 A 对角占优矛盾. 证毕. □

2. 非负矩阵和 M 矩阵

非负矩阵在迭代法中起着非常重要的作用. 一个非负矩阵就是其元素全为非负数的矩阵. 对两个 $m \times n$ 阶矩阵 A 和 B , 如果其所有相应的元素满足 $a_{ij} \geq b_{ij}$, 则称 $A \geq B$, 这样“ \geq ”就定义了矩阵集合上的一个偏序关系. 利用这种偏序关系, 非负矩阵 A 可以表示为 $A \geq O$. $A \geq B$ 也可表示为 $B \leq A$, 即“ \leq ”也定义了矩阵集合上的一个偏序关系.

性质 1.9 关于非负矩阵, 下列性质成立:

- (1) 若 $A \geq O, B \geq O$, 则 $AB \geq O, A+B \geq O$ 且 $A^k \geq O$.
- (2) 若 A, B, C 为非负矩阵且 $A \leq B$, 则 $AC \leq BC, CA \leq CB$.
- (3) 若 $O \leq A \leq B$, 则 $A^T \leq B^T$ 且 $A^k \leq B^k$ 对任意的 k 均成立.
- (4) 若 $O \leq A \leq B$, 则 $\|A\|_1 \leq \|B\|_1$, 且类似地 $\|A\|_\infty \leq \|B\|_\infty$.

定理 1.25 (Perron-Frobenius 定理) 设 $A \in \mathbb{R}^{n \times n}$, 则有下列结论成立:

- (1) 若 $A \geq O$, 则 A 有一个非负特征值等于 $\rho(A)$; 对 $\rho(A) > 0$, 相应的特征向量 $x \geq 0$, 且当 A 的任意元素增加时, $\rho(A)$ 不减少.
- (2) 若 $A \geq O$ 不可约, 则 A 有一个正特征值等于 $\rho(A)$; 对 $\rho(A) > 0$, 其对应的特征向量 $x > 0$; 当 A 的任意元素增加时, $\rho(A)$ 增加; 且 $\rho(A)$ 是 A 的一个简单特征值.

性质 1.10 设 $A \in \mathbb{R}^{n \times n}$, 则有下列性质成立:

- (1) 若 $O \leq A \leq B$, 则 $\rho(A) \leq \rho(B)$.
- (2) 若 $|A| \leq B$ ($|A|$ 表示 A 的元素取绝对值后得到的矩阵), 则 $\rho(A) \leq \rho(B)$.
- (3) 设 $A \geq O$, 则 $\rho(A) < 1$ 当且仅当 $I - A$ 非奇异且 $(I - A)^{-1} \geq O$.
- (4) 设 $A \geq O, x \geq 0$ 为非零向量, $\alpha > 0$. 若 $Ax > (\geq) \alpha x$, 则 $\rho(A) > (\geq) \alpha$; 若 $Ax < \alpha x$, 则 $\rho(A) < \alpha$.

定义 1.23 (1) 设 A 是一个 n 阶矩阵. 若 A 的主对角元是正的, 即 $a_{ii} > 0, i = 1, 2, \dots, n$, 而其非主对角元是非正的, 即 $a_{ij} \leq 0, i \neq j, i, j = 1, 2, \dots, n$, 则称 A 为 Z 矩阵. 若 Z 矩阵 A 为非奇异的, 且 $A^{-1} \geq O$, 则称 A 为 M 矩阵.

(2) 记矩阵 $\langle A \rangle = (\alpha_{ij})$, 其中主对角元 $\alpha_{ii} = |a_{ii}|$, 而非主对角元 $\alpha_{ij} = -|a_{ij}|$, 则称 $\langle A \rangle$ 为 A 的比较矩阵. 若 $\langle A \rangle$ 为非奇异的 M 矩阵, 则称 A 为 H 矩阵.

关于 M 矩阵和 H 矩阵, 有下面一些常用的性质.

性质 1.11 设 A 是 Z 矩阵, 则 A 是 M 矩阵当且仅当 $\rho(B) < 1$, 其中 $B = I - D^{-1}A$ 且 $D = \text{diag}(A)$.

性质 1.12 设 $A \in \mathbb{R}^{n \times n}$, 则有下列性质成立:

- (1) 严格对角占优或不可约对角占优矩阵是 H 矩阵.
- (2) 严格对角占优或不可约对角占优矩阵的 Z 矩阵是 M 矩阵.
- (3) 若 A 为 H 矩阵, 则 $|A|^{-1} \leq \langle A \rangle^{-1}$.
- (4) 若 $A = D - B$ 为 H 矩阵, $D = \text{diag}(A)$, 则 $|D|$ 非奇异, 且 $\rho(|D|^{-1}|B|) < 1$.

性质 1.13 设 $A, B \in \mathbb{R}^{n \times n}$ 是两个非奇异的 M 矩阵, 且 $A \leq B$, 则

(1) $A^{-1} \geq B^{-1}$.

(2) A 和 B 的每个主子矩阵也是 M 矩阵.

(3) 满足 $A \leq C \leq B$ 的 C 都是 M 矩阵. 特别地, 若 $A \leq C \leq \text{diag}(A)$, 则 C 是 M 矩阵.

性质 1.14 若 A 是 M 矩阵, B 是 Z 矩阵, 且 $A \leq B$, 则 B 也是 M 矩阵.

性质 1.15 (1) 若分块矩阵 $A = (A_{ij})_{i,j=1}^m$ 是 M 矩阵, 且其主对角块 A_{ii} 为方阵, 则 A_{ii} 也是 M 矩阵.

(2) 若 $A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$ 是 M 矩阵, 则 Schur 补 $S = A_{22} - A_{21}A_{11}^{-1}A_{12}$ 也是 M 矩阵.

3. 正定矩阵

关于正定矩阵, 在实矩阵的情形, 前面已作了介绍. 下面考虑复矩阵的情形.

定义 1.24 (1) 设矩阵 $A \in \mathbb{C}^{n \times n}$, 若对任意非零向量 $x \in \mathbb{C}^n$ 都有 $\text{Re}(x^H A x) > 0$, 则称 A 为正定的. 若 A 还是 Hermite 的, 则称为 Hermite 正定的, 记为 HPD.

(2) 设矩阵 $A \in \mathbb{C}^{n \times n}$, 若对任意非零向量 $x \in \mathbb{R}^n$ 都有 $x^T A x > 0$, 则称 A 为实正的.

任意 (实或复) 方阵 A 均可分裂为

$$A = H + iS, \quad (1.26)$$

式中:

$$H = \frac{1}{2}(A + A^H), \quad S = \frac{1}{2i}(A - A^H).$$

注意到 H 和 S 都是 Hermite 的, 而 iS 是反 Hermite 的. 矩阵 H 和 iS 分别称为 A 的 Hermite 部分和反 Hermite 部分. 当 A 是实矩阵且 u 是实向量时, 有 $(Au, u) = (Hu, u)$. 分裂 (1.26) 称为矩阵 A 的一个 Hermite-反 Hermite 分裂, 简称 HS 分裂.

正定矩阵具有下面一些常用的性质:

性质 1.16 关于矩阵 A , 下列性质成立:

(1) 设 A 为实正定矩阵, 则 A 是非奇异的, 且对任意的实向量 x , 存在标量 $\alpha > 0$ 满足 $(Ax, x) \geq \alpha \|x\|^2$.

(2) 正定矩阵 $A \in \mathbb{C}^{n \times n}$ 的所有特征特均有正的实部.

(3) 设 $A \in \mathbb{C}^{n \times n}$ 为正定矩阵, $P \in \mathbb{C}^{n \times n}$ 为 n 阶非奇异矩阵, 则 $P^H A P$ 亦为正定矩阵.

(4) 正定矩阵 $A \in \mathbb{C}^{n \times n}$ 的所有主子阵均为正定矩阵. 特别地, A 的所有对角元素均有正的实部.

(5) 设 $A \in \mathbb{C}^{n \times n}$, 则 A 是正定矩阵当且仅当 $A + A^H$ 为 Hermite 正定的.

定理 1.26 对于矩阵 A 的 HS 分裂 (1.26), 其特征值 λ_i 满足

$$\lambda_{\min}(H) \leq \operatorname{Re}(\lambda_i) \leq \lambda_{\max}(H), \quad \lambda_{\min}(S) \leq \operatorname{Im}(\lambda_i) \leq \lambda_{\max}(S).$$

当 B 为对称正定矩阵时, $\mathbb{C}^n \times \mathbb{C}^n$ 到 \mathbb{C} 上的映射 $x, y \rightarrow (x, y)_B \equiv (Bx, y)$ 是 \mathbb{C}^n 上的一个内积, 称此内积为能量内积. 而称其诱导范数为能量范数. 有时, 可以找到一个 HPD 矩阵 B , 使得一个给定的在欧几里得内积下的非 Hermite 矩阵 A 在能量内积下成为 Hermite 的, 即:

$$(Ax, y)_B = (x, Ay)_B, \quad \forall x, y.$$

最简单的例子是 $A = B^{-1}C$ 和 $A = CB$, 其中 C 是 Hermite 的且 B 是 Hermite 正定的.

1.6 模型问题: Poisson 问题

考虑二维规则区域上的 Poisson 问题

$$\begin{cases} -\left(\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2}\right) + c \frac{\partial u}{\partial x} + d \frac{\partial u}{\partial y} + eu = f(x, y), & (x, y) \in \Omega, \\ u(x, y)|_{\partial\Omega} = 0, \end{cases} \quad (1.27)$$

式中: $\Omega = (0, 1) \times (0, 1)$; f 为给定的函数; c, d, e 为非负常数.

设 x, y 方向均取 $n+1$ 个等距网格, 则内网格点共有 n^2 个, 步长为 $h = \frac{1}{n+1}$, 如图 1.1 所示.

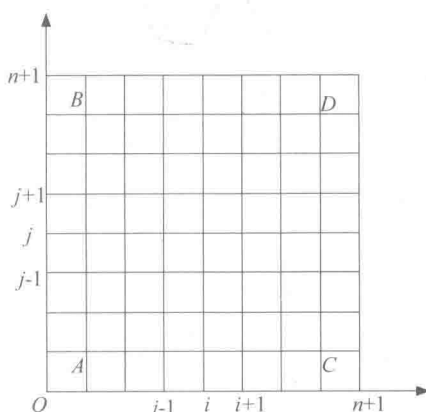


图 1.1 区域离散和节点标号

将

$$u_{xx}(x_i, y_j) \approx \frac{u_{i+1,j} - 2u_{i,j} + u_{i-1,j}}{h^2}, \quad u_{yy}(x_i, y_j) \approx \frac{u_{i,j+1} - 2u_{i,j} + u_{i,j-1}}{h^2},$$

$$u_x(x_i, y_j) \approx \frac{u_{i+1,j} - u_{i-1,j}}{2h}, \quad u_y(x_i, y_j) \approx \frac{u_{i,j+1} - u_{i,j-1}}{2h},$$

$$u(x_i, y_j) \approx u_{ij}, \quad f(x_i, y_j) \approx f_{ij}.$$

代入式 (1.27) 的第 1 式并整理, 得

$$(4 + h^2 e)u_{ij} - (u_{i+1,j} + u_{i-1,j} + u_{i,j+1} + u_{i,j-1}) + \frac{hc}{2}(u_{i+1,j} - u_{i-1,j}) \\ + \frac{hd}{2}(u_{i,j+1} - u_{i,j-1}) = h^2 f_{ij}, \quad i, j = 1, 2, \dots, n.$$

该方程组包含了 n^2 个未知量, 它的矩阵形式为

$$\mathbf{A}\mathbf{u} = \mathbf{f}, \quad (1.28)$$

式中: $\mathbf{A} = \mathbf{A}_1 + \mathbf{A}_2 + \mathbf{A}_3$ 为 n^2 阶的块三对角矩阵, 其具体形式为

$$\mathbf{A}_1 = \begin{bmatrix} \mathbf{B} & -\mathbf{I} & & & \\ -\mathbf{I} & \mathbf{B} & -\mathbf{I} & & \\ & \ddots & \ddots & \ddots & \\ & & -\mathbf{I} & \mathbf{B} & -\mathbf{I} \\ & & & -\mathbf{I} & \mathbf{B} \end{bmatrix}, \quad (1.29)$$

$$\mathbf{A}_2 = \frac{hc}{2} \begin{bmatrix} \mathbf{O} & \mathbf{I} & & & \\ -\mathbf{I} & \mathbf{O} & \mathbf{I} & & \\ & \ddots & \ddots & \ddots & \\ & & -\mathbf{I} & \mathbf{O} & \mathbf{I} \\ & & & -\mathbf{I} & \mathbf{O} \end{bmatrix}, \quad (1.30)$$

$$\mathbf{A}_3 = \frac{hd}{2} \begin{bmatrix} \mathbf{O} & -\mathbf{I} & & & \\ \mathbf{I} & \mathbf{O} & -\mathbf{I} & & \\ & \ddots & \ddots & \ddots & \\ & & \mathbf{I} & \mathbf{O} & -\mathbf{I} \\ & & & \mathbf{I} & \mathbf{O} \end{bmatrix}, \quad (1.31)$$

这里

$$\mathbf{B} = \begin{bmatrix} 4 + h^2 e & -1 & & & \\ -1 & 4 + h^2 e & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 4 + h^2 e & -1 \\ & & & -1 & 4 + h^2 e \end{bmatrix} \quad (1.32)$$

为对称三对角矩阵, \mathbf{I} 为 n 阶单位矩阵. n^2 维列向量 \mathbf{u} 定义为

$$\mathbf{u} = (u_{11}, u_{12}, \dots, u_{1n}, u_{21}, u_{22}, \dots, u_{2n}, \dots, u_{n1}, u_{n2}, \dots, u_{nn})^T, \quad (1.33)$$

网格中的每个内部节点都对应一个分量, u_{ij} 表示 $u(ih, jh)$ 的近似, 这些分量按沿 y 方向变化最快的方式排列. 图 1.1 中节点 A, B, C, D 的标号分别为 $(1, 1), (1, n), (n, 1), (n, n)$, 节点 (i, j) 对应的标号为 $i * n + j$, 它的上下左右的四个相邻节点为

$$i * n + (j + 1), \quad i * n + (j - 1), \quad (i - 1) * n + j, \quad (i + 1) * n + j.$$

右端向量 f 为

$$f = h^2(f_{11}, f_{12}, \cdots, f_{1n}, f_{21}, f_{22}, \cdots, f_{2n}, \cdots, f_{n1}, f_{n2}, \cdots, f_{nn})^T. \quad (1.34)$$

以后常常用此模型做数值实验, 故称它为模型问题.

习题 1

1.1 设 A 为正定矩阵, B 是与 A 同阶的 Hermite 矩阵. 试证明: $A + B$ 为正定矩阵的充分必要条件是 $A^{-1}B$ 的特征值都大于 1.

1.2 设矩阵 A 是可对角化的. 试证明: A 的特征值都是实数的充分必要条件是存在正定矩阵 B 使得 BA 为 Hermite 矩阵.

1.3 设 A 和 B 为同阶半正定矩阵. 证明: $\det(A + B) \geq \det(A) + \det(B)$.

1.4 设 A 为反对称矩阵. 证明:

- (1) 矩阵 $I - A$ 是非奇异的;
- (2) 矩阵 $C = (I - A)^{-1}(I + A)$ 是正交矩阵.

1.5 设 A 和 B 为同阶幂等矩阵, 并且满足 $\mathcal{R}(A) = \mathcal{R}(B)$ 和 $\mathcal{N}(A) = \mathcal{N}(B)$. 证明: $A = B$.

1.6 设

$$A = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix},$$

求矩阵 A 的奇异值分解.

1.7 设矩阵 $A \in \mathbb{C}^{m \times n}$ 的一个满秩分解为 $A = FG$, 求矩阵 $\begin{bmatrix} A & A \\ A & A \end{bmatrix}$ 的一个满秩分解.

1.8 设非零向量 $x \in \mathbb{C}^n$. 证明:

- (1) $\|x\|_\infty \leq \|x\|_1 \leq n\|x\|_\infty$; (2) $\frac{1}{\sqrt{n}}\|x\|_1 \leq \|x\|_2 \leq \|x\|_1$;
- (3) $\|x\|_\infty \leq \|x\|_2 \leq \sqrt{n}\|x\|_\infty$.

1.9 设 A 为 n 阶矩阵. 证明:

- (1) $\frac{1}{n}\|A\|_\infty \leq \|A\|_1 \leq n\|A\|_\infty$; (2) $\frac{1}{\sqrt{n}}\|A\|_2 \leq \|A\|_1 \leq \sqrt{n}\|A\|_2$;
- (3) $\frac{1}{\sqrt{n}}\|A\|_2 \leq \|A\|_\infty \leq \sqrt{n}\|A\|_2$; (4) $\frac{1}{\sqrt{n}}\|A\|_F \leq \|A\|_2 \leq \|A\|_F$.

1.10 设 $\mathbf{0} \neq \mathbf{v} \in \mathbb{R}^n$, $\mathbf{A} \in \mathbb{R}^{n \times n}$. 证明:

$$\|\mathbf{A}(\mathbf{I} - (\mathbf{v}^T \mathbf{v})^{-1} \mathbf{v} \mathbf{v}^T)\|_F^2 = \|\mathbf{A}\|_F^2 - \frac{\|\mathbf{A} \mathbf{v}\|_2^2}{\mathbf{v}^T \mathbf{v}}.$$

1.11 证明: $\rho(\mathbf{A}) < 1$ 当且仅当存在正定矩阵 \mathbf{B} , 使得 $\mathbf{B} - \mathbf{A} \mathbf{B} \mathbf{A}^H$ 也为正定矩阵.

1.12 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 且 $\mathbf{A} \geq \mathbf{O}$.

(1) 若 \mathbf{A} 的各行元素之和相等, 则 $\rho(\mathbf{A}) = \|\mathbf{A}\|_\infty$;

(2) 若 \mathbf{A} 的各列元素之和相等, 则 $\rho(\mathbf{A}) = \|\mathbf{A}\|_1$.

1.13 设 \mathbf{A} 为 Z 矩阵. 证明: \mathbf{A} 是非奇异的 M 矩阵当且仅当存在 $\mathbf{x} > \mathbf{0}$ 使得 $\mathbf{A} \mathbf{x} > \mathbf{0}$.

1.14 设 \mathbf{A} 是 Hermite 矩阵, 且 \mathbf{X} 是 \mathbf{A} 的广义逆. 证明: \mathbf{X}^2 是 \mathbf{A}^2 的广义逆.

1.15 证明: $\mathbf{B}^\dagger \mathbf{A}^\dagger = (\mathbf{A} \mathbf{B})^\dagger$ 当且仅当 $\mathcal{R}(\mathbf{A}^H \mathbf{A} \mathbf{B}) \subseteq \mathcal{R}(\mathbf{B})$, 且 $\mathcal{R}(\mathbf{B} \mathbf{B}^H \mathbf{A}^H) \subseteq \mathcal{R}(\mathbf{A}^H)$.

1.16 试证明下述结论:

(1) $\mathbf{A}^\dagger = (\mathbf{A}^H \mathbf{A})^\dagger \mathbf{A}^H = \mathbf{A}^H (\mathbf{A} \mathbf{A}^H)^\dagger$;

(2) 设 $\mathbf{a}, \mathbf{b} \in \mathbb{C}^n$, 则 $(\mathbf{a} \mathbf{b}^H)^\dagger = (\mathbf{a}^H \mathbf{a})^\dagger (\mathbf{b}^H \mathbf{b})^\dagger \mathbf{b} \mathbf{a}^H$.

1.17 设 \mathbf{A} 为正规矩阵. 证明: $\mathbf{A} \mathbf{A}^\dagger = \mathbf{A}^\dagger \mathbf{A}$, 且对任一自然数 k , 有 $(\mathbf{A}^k)^\dagger = (\mathbf{A}^\dagger)^k$.

1.18 设 \mathbf{A} 为 Z 矩阵. 证明: \mathbf{A} 是 M 矩阵当且仅当 $\rho(\mathbf{B}) < 1$, 其中 $\mathbf{B} = \mathbf{I} - \mathbf{D}^{-1} \mathbf{A}$ 且 $\mathbf{D} = \text{diag}(\mathbf{A})$.

第 2 章 正交变换和投影方法

向量序列的正交化和矩阵的正交 (相似) 变换在数值代数中占有独特的地位. QR 分解可用于线性方程组特别是最小二乘问题的求解. 而基于 QR 分解的 QR 方法可用于计算矩阵的特征值和特征向量. 此外, Krylov 子空间的正交化可用于构造线性方程组的迭代法. 投影方法的思想则是构造线性方程组迭代法和矩阵特征值计算的基础和出发点. 本章将主要介绍两种常用的正交变换、QR 分解、线性无关向量组的正交化以及 Krylov 子空间的性质及其正交化, 最后介绍投影方法的一般框架.

2.1 两种常用的正交变换

2.1.1 Householder 变换

定义 2.1 设 $u \in \mathbb{R}^n$ 满足 $\|u\|_2 = 1$, 称 n 阶矩阵

$$H = I - 2uu^T \quad (2.1)$$

为 Householder 矩阵 (初等反射矩阵), u 称为 Householder 向量.

下面的定理给出了 Householder 矩阵的性质.

定理 2.1 设 H 为式 (2.1) 定义的 Householder 矩阵, 则

- (1) $\det(H) = -1$.
- (2) $H^T = H$, $H^T H = I$, $H^{-1} = H$, $H^2 = I$.
- (3) H 仅有两个互不相同的特征值 -1 和 1 , 且 -1 是单重的, 相应的特征向量为 u . 而 1 是 $n-1$ 重的, 相应的特征向量为所有与 u 正交的非零向量.
- (4) 设 $x, u \in \mathbb{R}^n$ 满足 $u^T x = 0$, $\alpha \in \mathbb{R}$. 则

$$H(x + \alpha u) = x - \alpha u.$$

证明 (1) 利用分块矩阵的乘法运算可以方便地验证 $\det(H) = -1$. 事实上, 对

$$\begin{aligned} & \begin{bmatrix} I & 0 \\ -u^T & 1 \end{bmatrix} \begin{bmatrix} I & 2u \\ 0 & 1 \end{bmatrix} \begin{bmatrix} I - 2uu^T & 0 \\ u^T & 1 \end{bmatrix} \\ &= \begin{bmatrix} I & 0 \\ -u^T & 1 \end{bmatrix} \begin{bmatrix} I & 2u \\ u^T & 1 \end{bmatrix} = \begin{bmatrix} I & 2u \\ 0 & 1 - 2u^T u \end{bmatrix} \end{aligned}$$

两端取行列式, 得 $\det(H) = 1 - 2u^T u = -1$.

(2) 利用 H 的定义, 容易验证这四个等式.

(3) 由于 $Hu = (I - 2uu^T)u = -u$, 故 -1 是 H 的一个特征值, 且几何重数至少为 1, u 为相应的特征向量. 另外, 对于任意与 u 正交的 n 维向量 x , 有 $Hx = x$, 即 1 为 H 的一个特征值, 且几何重数至少为 $n-1$, 与 u 正交的任一非零向量作为相应

的特征向量. 由于特征值的几何重数不超过代数重数, 故特征值 -1 与 1 的代数重数分别至少为 1 与 $n-1$, 其和至少为 n . 注意到矩阵的特征值代数重数之和不会超过 n , 故 -1 与 1 的代数重数分别刚好为 1 与 $n-1$.

(4) 直接计算

$$H(x + \alpha u) = (I - 2uu^T)(x + \alpha u) = x - \alpha u.$$

证毕. □

下面的定理是 Householder 矩阵的另一个性质.

定理 2.2 设 H_{n-m} 是 $n-m$ 阶 Householder 矩阵, 则

$$H = \begin{bmatrix} I_m & O \\ O & H_{n-m} \end{bmatrix}$$

是 n 阶 Householder 矩阵.

证明 设 u_{n-m} 是 $n-m$ 维单位列向量, 则有

$$H_{n-m} = I_{n-m} - 2u_{n-m}u_{n-m}^T,$$

$$\begin{aligned} H &= \begin{bmatrix} I_m & O \\ O & I_{n-m} - 2u_{n-m}u_{n-m}^T \end{bmatrix} \\ &= \begin{bmatrix} I_m & O \\ O & I_{n-m} \end{bmatrix} - 2 \begin{bmatrix} O & O \\ O & u_{n-m}u_{n-m}^T \end{bmatrix} \\ &= \begin{bmatrix} I_m & O \\ O & I_{n-m} \end{bmatrix} - 2 \begin{bmatrix} 0 \\ u_{n-m} \end{bmatrix} \begin{bmatrix} 0 & u_{n-m}^T \end{bmatrix} \\ &= I_n - 2u_n u_n^T, \end{aligned}$$

式中: $u_n = \begin{bmatrix} 0 \\ u_{n-m} \end{bmatrix} \in \mathbb{R}^n.$

由于 $u_n^T u_n = u_{n-m}^T u_{n-m} = 1$, 所以 H 是 n 阶 Householder 矩阵. 证毕. □

定理 2.3 设 $x = (x_1, x_2, \dots, x_n)^T \in \mathbb{R}^n$ ($n > 1$) 为任意的非零向量, 则存在 $u \in \mathbb{R}^n$ 满足 $\|u\|_2 = 1$, 使得式 (2.1) 定义的 Householder 矩阵满足

$$Hx = \alpha e_1, \tag{2.2}$$

式中: $\alpha = \pm \|x\|_2$. 且使得式 (2.2) 成立的 u 在相差一个符号的意义下是唯一确定的.

证明 由定义 2.1, $H = I - 2uu^T$, 于是

$$Hx = x - 2(u^T x)u.$$

欲使式 (2.2) 成立, 必须满足

$$2(u^T x)u = x - \alpha e_1.$$

由于 $\|u\|_2 = 1$, 可取

$$u = \frac{x - \alpha e_1}{\|x - \alpha e_1\|_2}.$$

又因 H 是正交矩阵, 为使式 (2.2) 成立, 必须有

$$\|x\|_2 = \|Hx\|_2 = \|\alpha e_1\|_2 = |\alpha| \cdot \|e_1\|_2 = |\alpha|,$$

即 $\alpha = \pm \|x\|_2$. 容易验证, 如上选取的 H 满足式 (2.2). 事实上, 有

$$\begin{aligned} Hx &= x - 2 \frac{(x - \alpha e_1)(x - \alpha e_1)^T}{\|x - \alpha e_1\|_2^2} x \\ &= x - \frac{2(x - \alpha e_1)^T x}{\|x - \alpha e_1\|_2^2} (x - \alpha e_1) \\ &= x - \frac{2\|x\|_2^2 - 2\alpha e_1^T x}{\|x\|_2^2 - 2\alpha e_1^T x + \alpha^2} (x - \alpha e_1) \\ &= x - \frac{2\alpha^2 - 2\alpha e_1^T x}{\alpha^2 - 2\alpha e_1^T x + \alpha^2} (x - \alpha e_1) \\ &= x - (x - \alpha e_1) = \alpha e_1. \end{aligned}$$

此外, 由 u 的选取过程知, 这样的 u 在相差一个符号的意义下是唯一确定的. 证毕. \square

注 2.1 由定理 2.3, 可按如下步骤来构造确定 H 的单位向量 u :

- (1) 计算 $v = x \pm \|x\|_2 e_1$.
- (2) 计算 $u = v / \|v\|_2$.

上述计算涉及 $\|x\|_2$ 前的符号选取问题. 如果选取

$$v = x - \|x\|_2 e_1,$$

就会出现计算

$$v_1 = x_1 - \|x\|_2$$

的问题, 其中 v_1, x_1 分别表示 v, x 的第 1 个分量. 在 $x_1 > 0$ 时, 按上式计算 v_1 可能会导致有效数字的丢失. 在这种情形可改用下面的计算方式, 即

$$v_1 = x_1 \pm \|x\|_2 = \frac{x_1^2 - \|x\|_2^2}{x_1 + \|x\|_2} = \frac{-(x_2^2 + \cdots + x_n^2)}{x_1 + \|x\|_2},$$

并且只要在 $x_1 > 0$ 时使用这一公式计算, 就可以避免两个相近的数相减的情形.

注意到

$$H = I - 2uu^T = I - \frac{2}{v^T v} vv^T = I - \beta vv^T,$$

式中: $\beta = 2/(v^T v)$. 这样就没有必要显式地求出向量 u , 只需求出 β 和 v 即可. 在实际计算中, 往往将 v 归化为第 1 个分量为 1 的向量 (只需作变换 $v := v/v_1$ 即可). 这样做的优点是可以把 v 的后 $n-1$ 个分量保存在 x 的后 $n-1$ 个化为 0 的分量位置上, 而 v 的第一个分量 1 就无需保存了.

在计算时, 溢出现象也是必须要考虑的问题. 在上述计算中, 如果 x 某分量过大, 其平方运算可能会出现上溢的现象. 为了解决这个问题, 可以用 $x/\|x\|_\infty$ 代替 x 来构造 v , 相当于在原来的 v 之前乘了一个常数 $\alpha = 1/\|x\|_\infty$, 而 αv 与 v 的单位化向量是相同的.

基于上述讨论, 可得如下基本算法.

算法 2.1 本算法计算实 Householder 矩阵 $H = I - \beta vv^T$ 中满足 $v_1 = 1$ 的 v 和 β .

- 步 1, 输入 n 维实向量 x . 计算 $\eta = \|x\|_\infty$, 置 $x := x/\eta$.
- 步 2, $v := [1, x(2:n)]^T$, 计算 $\sigma = x_2^2 + \cdots + x_n^2$.
- 步 3, 对于 $\sigma = 0$, 若 $x_1 \geq 0$, $\beta := 0$, 否则 $\beta := 2$, 终止计算.
- 步 4, 对于 $\sigma > 0$, 计算 $\alpha := \sqrt{x_1^2 + \sigma}$. 若 $x_1 \leq 0$, $v_1 = x_1 - \alpha$, 否则 $v_1 = -\sigma/(x_1 + \alpha)$.
- 步 5, 计算 $\beta := 2v_1^2/(\sigma + v_1^2)$, $v := v/v_1$.

执行算法 2.1 的运算量约为 $4n$. 进一步, 根据上述算法可编制 MATLAB 程序如下:

```
%实向量Householder变换程序-r_house.m
function [v,beta]=r_house(x)
%本函数计算Householder矩阵H=I-beta*v*v'中满足v(1)=1的v和beta.
n=length(x);
eta=norm(x,inf); x=x/eta;
sigma=x(2:n)'*x(2:n);
v=[1; x(2:n)];
if sigma==0
    if x(1)>=0
        beta=0;
    else
        beta=2;
    end
else
    alpha=(x(1)^2+sigma)^0.5;
    if x(1)<=0
        v(1)=x(1)-alpha;
```

```

else
    v(1)=-sigma/(x(1)+alpha);
end
beta=2*v(1)^2/(sigma+v(1)^2);
v=v/v(1);
end

```

例 2.1 已知向量 $x = (2, 5, 7, 1)^T$, 构造 Householder 矩阵 H 使得 $Hx = \|x\|_2 e_1$.

解 在 MATLAB 命令窗口依次运行如下代码:

```

>> x=[2,5,7,1]'; [v,beta]=r_house(x);
>> H=eye(length(x))-beta*v*v';
>> y=H*x

```

即得所需求的结果.

注 2.2 在应用 Householder 变换约化一个给定的矩阵为某一需要的形式时, 利用其特殊结构是非常重要的. 当计算 Householder 矩阵 $H = I - \beta vv^T \in \mathbb{R}^{m \times m}$ 与一个已知矩阵 $A \in \mathbb{R}^{m \times n}$ 的乘积时, 实际计算中 H 并不需要以显式的方式给出, 而是根据如下的公式来计算:

$$HA = (I - \beta vv^T)A = A - \beta v(A^T v)^T = A - vw^T,$$

式中: $w = \beta A^T v$. 换言之, 可以按下列两个步骤计算 HA :

(1) 计算 $w = \beta A^T v$.

(2) 计算 $B = A - vw^T$.

矩阵 B 即为所求的乘积 HA . 完成这一计算任务所需要的运算量为 $4mn$.

注 2.3 在复数域情形, Householder 变换是指如下形式的矩阵

$$H(w) = I - 2ww^H, \quad (2.3)$$

式中: $w \in \mathbb{C}^n$ 满足 $w^H w = 1$. 相应于实 Householder 矩阵是正交矩阵, 复 Householder 矩阵是一个酉矩阵. 事实上, 复 Householder 变换除具备实 Householder 变换相应的 (1) ~ (4) 条性质外, 还有如下性质:

(5) 设 $x, y \in \mathbb{C}^n$, $x \neq y$. 则存在 Householder 变换 $H(w)$ 使得 $H(w)x = y$ 的充分必要条件是

$$x^H x = y^H y, \quad x^H y = y^H x.$$

并且在上述条件成立时, 所需的向量 w 可取为

$$w = e^{i\theta} \frac{x - y}{\|x - y\|_2},$$

式中: θ 为任意实数.

由注 2.3, 若取 $\mathbf{x} = (x_1, x_2, \dots, x_n)^T \neq \mathbf{0}$, $\mathbf{y} = -\bar{\tau}\|\mathbf{x}\|_2\mathbf{e}_1$, 其中

$$\tau = \begin{cases} 1, & \text{当 } x_1 = 0 \text{ 时,} \\ \frac{\bar{x}_1}{|x_1|}, & \text{当 } x_1 \neq 0 \text{ 时,} \end{cases} \quad (2.4)$$

则满足 $\mathbf{x}^H\mathbf{x} = \mathbf{y}^H\mathbf{y}$, $\mathbf{x}^H\mathbf{y} = \mathbf{y}^H\mathbf{x}$. 此时成立

$$\mathbf{H}(\mathbf{w})\mathbf{x} = -\bar{\tau}\|\mathbf{x}\|_2\mathbf{e}_1,$$

式中: \mathbf{w} 可取为

$$\mathbf{w} = \frac{\tau\mathbf{x} + \gamma\mathbf{e}_1}{\|\tau\mathbf{x} + \gamma\mathbf{e}_1\|_2}, \quad \gamma = \|\mathbf{x}\|_2.$$

而且易见当 \mathbf{x} 为实向量时, 这样得到的 $\mathbf{H}(\mathbf{w})$ 是实对称正交矩阵.

基于以上讨论, 有下面的算法.

算法 2.2 给定一个向量 $\mathbf{x} \in \mathbb{C}^n$, 本算法计算一个向量 \mathbf{w} 和一个数 γ , 使得 $\mathbf{H}\mathbf{x} = \gamma\mathbf{e}_1$, 其中 $\mathbf{H} = \mathbf{I} - \mathbf{w}\mathbf{w}^H$, $\|\mathbf{w}\|_2 = \sqrt{2}$.

function $[\mathbf{w}, \gamma] = \mathbf{c_house}(\mathbf{x})$

$\mathbf{w} = \mathbf{x}; \gamma = \|\mathbf{x}\|_2;$

if $\gamma = 0$

$w(1) = \sqrt{2};$ 结束

end

if $w(1) \neq 0$

$\tau = \overline{w(1)}/|w(1)|;$

else

$\tau = 1;$

end

$\mathbf{w} = (\tau/\gamma)\mathbf{w}; w(1) = w(1) + 1;$

$\mathbf{w} = \mathbf{w}/\sqrt{w(1)}; \gamma = -\bar{\tau}\gamma;$

注 2.4 算法 2.2 是用类似于 MATLAB 语句给出的, 后文的许多算法也将采用这种形式给出. 因此, 熟悉 MATLAB 语言是很重要的. 可以说, 用类似于 MATLAB 语言的语句来描述算法往往要比用文字和公式描述算法精到得多.

注 2.5 算法 2.1 中的 Householder 变换只是针对实向量的, 而算法 2.2 既能用于实向量又能用于复向量. 这一点是很重要的, 因为许多实际问题, 如时谐涡流场的计算, 其离散形式是一个大型稀疏且具有特殊块结构的复线性方程组. 此时使用 Krylov 子空间方法 (如 GMRES 方法) 求解时, 需要用到针对复向量的 Householder 变换.

算法 2.2 的 MATLAB 程序如下:


```

%复向量的Householder变换程序-c_house.m
function [w,gamma]=c_house(x)
%给定复向量x, 本函数计算一个向量w和一个数gamma
%满足H*x=gamma*e1, 其中H=I-w*w', ||w||=sqrt(2).
w=x; gamma=norm(x);
if gamma==0
    w(1)=sqrt(2); return;
end
if w(1)==0
    tau=1;
else
    tau=conj(w(1))/abs(w(1));
end
w=(tau/gamma)*w; w(1)=w(1)+1;
w=w/sqrt(w(1)); gamma=-conj(tau)*gamma;

```

例 2.2 已知复向量 $x = (2 + i, 5 - 3i, 7, 1 + 2i)^T$, 构造 Householder 矩阵 H 使得 $Hx = \|x\|_2 e_1$.

解 在 MATLAB 命令窗口依次运行如下代码:

```

>> x=[2+i,5-3*i,7,1+2*i]'; [w,gama]=c_house(x);
>> H=eye(length(x))-w*w'; y=H*x

```

即得所需求的结果.

2.1.2 Givens 变换

Householder 变换可以将一个向量中若干相邻分量约化为 0. 但如果将向量中的某一个分量化为 0, 则采用 Givens 变换更为有效. Givens 变换 (矩阵) 定义如下.

定义 2.2 设实数 c 和 s 满足 $c^2 + s^2 = 1$, 称 n 阶矩阵

$$G_{ik}(c, s) = \begin{bmatrix} 1 & & & & & & \\ & \ddots & & & & & \\ & & c & & s & & \\ & & & \ddots & & & \\ & & & & c & & \\ & & -s & & c & & \\ & & & & & \ddots & \\ & & & & & & 1 \end{bmatrix} \quad \begin{matrix} (i) \\ (k) \end{matrix} \quad (i \neq k) \quad (2.5)$$

为 Givens 变换 (矩阵), 或初等旋转矩阵.

容易证明, Givens 矩阵具有如下性质.

定理 2.4 设 G 是由式 (2.5) 定义的 Givens 矩阵. 则

(1) $G^T G = I$, 即 $G_{ik}(c, s)$ 为正交矩阵.

(2) $\det(G) = 1$.

(3) $G_{ik}(c, s)^{-1} = G_{ik}(c, s)^T = G_{ik}(c, -s)$.

(4) 对于任意的 $x \in \mathbb{R}^n$, Givens 变换 $y = G_{ik}(c, s)x$ 只改变 x 的第 i, k 个分量.

证明 结论 (1) ~ (3) 容易验证. 只证明结论 (4). 设 $x = (x_1, x_2, \dots, x_n)^T$, $y = (y_1, y_2, \dots, y_n)^T$, 则由 $y = G_{ik}(c, s)x$, 得

$$y_i = cx_i + sx_k, \quad y_k = -sx_i + cx_k, \quad y_j = x_j, \quad (j \neq i, k).$$

证毕. □

由定理 2.4 的结论 (4) 不难发现, 当 $x_i^2 + x_k^2 \neq 0$ 时, 选取

$$c = \frac{x_i}{\sqrt{x_i^2 + x_k^2}}, \quad s = \frac{x_k}{\sqrt{x_i^2 + x_k^2}}, \quad (2.6)$$

可使 $y_i = \sqrt{x_i^2 + x_k^2} > 0$, $y_k = 0$. 由此, 给定 x 第 i 个和第 k 个分量, 可按下列步骤计算 c 和 s 使得 $G_{ik}(c, s)x$ 的第 k 个分量为 0.

算法 2.3 (计算实 Givens 变换) 给定两个实数 a 和 b , 本算法计算三个数 c, s, η , 使得

$$\begin{bmatrix} c & s \\ -s & c \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} \eta \\ 0 \end{bmatrix}.$$

function $[c, s, \eta] = \text{r_givens}(a, b)$

if $b = 0$

$c = 1; s = 0; \eta = a$; 结束

end

if $a = 0$

$c = 0; s = 1; \eta = b$; 结束

end

if $|b| \geq |a|$

$t = a/b; s = 1/\sqrt{1+t^2}; c = st; \eta = |b|/s;$

else

$t = b/a; c = 1/\sqrt{1+t^2}; s = ct; \eta = |a|/c;$

end

执行算法 2.3 只需 5 次乘除法、1 次加法和 1 次开平方运算. 如果不计算 η , 则可以减少一次除法运算. 进一步, 根据上述算法可编制 MATLAB 程序如下:

%Given变换程序-r_givens.m

function $[c, s, \eta] = \text{r_givens}(a, b)$

%计算 c, s 满足 $\begin{bmatrix} c & s \\ -s & c \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} \eta \\ 0 \end{bmatrix}$

```

if b==0, c=1; s=0; eta=a; end
if a==0, c=0; s=1; eta=b; end
if abs(b)>abs(a)
    t=a/b; s=1/sqrt(1+t^2); c=s*t; eta=abs(b)/s;
else
    t=b/a; c=1/sqrt(1+t^2); s=c*t; eta=abs(a)/c;
end

```

下面定理的结论表明 Givens 变换具有与 Householder 变换相同的功能.

定理 2.5 设 $\mathbf{x} = (x_1, x_2, \dots, x_n)^T \neq \mathbf{0}$, 则存在有限个 Givens 矩阵的乘积 \mathbf{G} , 使得 $\mathbf{G}\mathbf{x} = \|\mathbf{x}\|_2 \mathbf{e}_1$.

证明 (1) 设 $x_1 \neq 0$, 依次构造

$$\mathbf{G}_{12}: c = \frac{x_1}{\sqrt{x_1^2 + x_2^2}}, \quad s = \frac{x_2}{\sqrt{x_1^2 + x_2^2}},$$

$$\Rightarrow \mathbf{G}_{12}\mathbf{x} = \left(\sqrt{x_1^2 + x_2^2}, 0, x_3, \dots, x_n \right)^T.$$

$$\mathbf{G}_{13}: c = \frac{\sqrt{x_1^2 + x_2^2}}{\sqrt{x_1^2 + x_2^2 + x_3^2}}, \quad s = \frac{x_3}{\sqrt{x_1^2 + x_2^2 + x_3^2}},$$

$$\Rightarrow \mathbf{G}_{13}(\mathbf{G}_{12}\mathbf{x}) = \left(\sqrt{x_1^2 + x_2^2 + x_3^2}, 0, 0, x_4, \dots, x_n \right)^T.$$

⋮

$$\mathbf{G}_{1n}: c = \frac{\sqrt{x_1^2 + \dots + x_{n-1}^2}}{\sqrt{x_1^2 + \dots + x_n^2}}, \quad s = \frac{x_n}{\sqrt{x_1^2 + \dots + x_n^2}},$$

$$\Rightarrow \mathbf{G}_{1n}(\mathbf{G}_{1,n-1} \dots \mathbf{G}_{13}\mathbf{G}_{12}\mathbf{x}) = \left(\sqrt{x_1^2 + \dots + x_n^2}, 0, \dots, 0 \right)^T.$$

令 $\mathbf{G} = \mathbf{G}_{1n}\mathbf{G}_{1,n-1} \dots \mathbf{G}_{13}\mathbf{G}_{12}$, 则有 $\mathbf{G}\mathbf{x} = \|\mathbf{x}\|_2 \mathbf{e}_1$.

(2) 设 $x_1 = \dots = x_{i-1} = 0, x_i \neq 0 (1 < i \leq n)$, 则由 \mathbf{G}_{1i} 开始即可. 证毕. □

与 Householder 矩阵类似, 用 Givens 矩阵左乘或右乘一个已知矩阵 \mathbf{A} , 利用其特殊结构也是极为有利的. 假定 $\mathbf{A} \in \mathbb{R}^{m \times n}$. 如果 $\mathbf{G}_{ik}(c, s) \in \mathbb{R}^{m \times m}$, 则用 $\mathbf{G}_{ik}(c, s)\mathbf{A}$ 修正 \mathbf{A} 仅影响 \mathbf{A} 的 i, k 两行, 可以用

$$\begin{bmatrix} a_{i1} & a_{i2} & \cdots & a_{in} \\ a_{k1} & a_{k2} & \cdots & a_{kn} \end{bmatrix} := \begin{bmatrix} c & s \\ -s & c \end{bmatrix} \begin{bmatrix} a_{i1} & a_{i2} & \cdots & a_{in} \\ a_{k1} & a_{k2} & \cdots & a_{kn} \end{bmatrix}$$

实现, 这仅需 $6n$ 个 flop (浮点运算). 实现上述运算的 MATLAB 程序段如下:

```
for j=1:n
```

```

t1=A(i,j); t2=A(k,j);
A(i,j)=c*t1+s*t2;
A(k,j)=-s*t1+c*t2;
end

```

同样, 如果 $G_{ik}(c, s) \in \mathbb{R}^{n \times n}$, 则用 $AG_{ik}(c, s)$ 修正 A 仅影响 A 的 i, k 两列, 可以用

$$\begin{bmatrix} a_{1i} & a_{1k} \\ a_{2i} & a_{2k} \\ \vdots & \vdots \\ a_{mi} & a_{mk} \end{bmatrix} := \begin{bmatrix} a_{1i} & a_{1k} \\ a_{2i} & a_{2k} \\ \vdots & \vdots \\ a_{mi} & a_{mk} \end{bmatrix} \begin{bmatrix} c & s \\ -s & c \end{bmatrix}$$

实现, 这只需 $6m$ 个 flop. 实现上述运算的 MATLAB 程序段如下:

```

for j=1:m
    t1=A(j,i); t2=A(j,k);
    A(j,i)=c*t1-s*t2;
    A(j,k)=s*t1+c*t2;
end

```

注 2.6 复数情形的 Givens 矩阵是指定义 2.2 中的 c 和 s 为复数且满足 $|c|^2 + |s|^2 = 1$. 此时的 Givens 矩阵 $G_{ik}(c, s)$ 是一个酉矩阵. 设 $x = (x_1, x_2, \dots, x_n)^T \in \mathbb{C}^n$, $y = G_{ik}(c, s)x = (y_1, y_2, \dots, y_n)^T$. 容易验证, 当 $|x_i|^2 + |x_k|^2 \neq 0$ 时, 若取

$$c = \frac{|x_i|}{\sqrt{|x_i|^2 + |x_k|^2}}, \quad s = \frac{\bar{x}_k}{\sqrt{|x_i|^2 + |x_k|^2}} \frac{x_i}{|x_i|}, \quad (2.7)$$

则有

$$y_s = x_s, \quad s \neq i, k, \quad y_i = \frac{x_i}{|x_i|} \sqrt{|x_i|^2 + |x_k|^2}, \quad y_k = 0, \quad (2.8)$$

即可以通过复 Givens 变换将复向量 x 的第 k 个分量化为零.

根据式 (2.7), 可设计算法如下.

算法 2.4 (计算复 Givens 变换) 给定两个复数 a 和 b , 本算法计算三个数 c, s, η , 使得

$$\begin{bmatrix} c & s \\ -\bar{s} & \bar{c} \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} \eta \\ 0 \end{bmatrix}.$$

function $[c, s, \eta] = \text{c.givens}(a, b)$

if $b = 0$

$c = 1; s = 0; \eta = a;$ 结束

end

if $a = 0$

$c = 0; s = 1; \eta = b;$ 结束

end

$u = a/|a|$; $t = \sqrt{|a|^2 + |b|^2}$;

$c = |a|/t$; $s = ub/t$; $\eta = ut$;

算法 2.4 的 MATLAB 程序如下:

%复Given变换程序-c_givens.m

function [c,s,eta]=c_givens(a,b)

%计算c,s满足 $[c,s;-s',c']*[a;b]=[eta;0]$

if b==0, c=1; s=0; eta=a; end

if a==0, c=0; s=1; eta=b; end

u=a/abs(a); t=sqrt(abs(a)^2+abs(b)^2);

c=abs(a)/t; s=u*conj(b)/t; eta=u*t;

2.2 QR 分解

实现矩阵 QR 分解最常用的方法有三种, 分别是 Householder 正交化方法、Givens 正交化方法和 Gram-Schmidt 正交化方法. 本节先给出前两种方法的具体实现.

定义 2.3 设 $A \in \mathbb{R}^{m \times n}$ ($m \geq n$) 为列满秩矩阵, 若有正交矩阵 $Q \in \mathbb{R}^{m \times m}$ 与上三角矩阵 $R \in \mathbb{R}^{m \times n}$ 使得 $A = QR$, 则称 QR 为 A 的 QR 分解.

2.2.1 Householder 变换 QR 分解

本节介绍如何使用 Householder 变换求矩阵的 QR 分解.

定理 2.6 设 $A \in \mathbb{R}^{m \times n}$ ($m \geq n$) 为列满秩矩阵, 则存在有限个 Householder 矩阵的乘积 Q , 使得

$$A = Q\tilde{R} = Q \begin{bmatrix} R \\ O \end{bmatrix}, \quad (2.9)$$

式中: $R \in \mathbb{R}^{n \times n}$ 为非奇异的上三角矩阵.

证明 取 $A_1 := A$. 因为 $A = (a_{ij}) \in \mathbb{R}^{m \times n}$ 列满秩, 故 A 的第 1 列 $a_1 \neq 0$. 于是可以构造 Householder 矩阵 $H_1 \in \mathbb{R}^{m \times m}$, 使得 $H_1 a_1 = \alpha_1 e_1$, 这里 $\alpha_1 = \|a_1\|_2$, $e_1 \in \mathbb{R}^m$. 则有

$$A_2 = H_1 A_1 = \begin{bmatrix} \alpha_1 & * \\ 0 & A_{22}^{(2)} \end{bmatrix}.$$

易见 $A_{22}^{(2)} \in \mathbb{R}^{(m-1) \times (n-1)}$ 列满秩. 对 $A_{22}^{(2)}$ 第 1 列的 $m-1$ 维非零向量 a_2 构造 $m-1$ 阶 Householder 矩阵 \tilde{H}_2 , 使得 $\tilde{H}_2 a_2 = \alpha_2 e_1$, 这里 $\alpha_2 = \|a_2\|_2$, $e_1 \in \mathbb{R}^{m-1}$. 令

$H_2 = \text{diag}(1, \widetilde{H}_2)$, 则有

$$A_3 = H_2 A_2 = \begin{bmatrix} \alpha_1 & * & * & \cdots & * \\ 0 & \alpha_2 & * & \cdots & * \\ 0 & 0 & * & \cdots & * \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & * & \cdots & * \end{bmatrix} = \begin{bmatrix} A_{11}^{(3)} & A_{12}^{(3)} \\ O & A_{22}^{(3)} \end{bmatrix} \begin{matrix} 2 \\ m-2 \\ 2 & n-2 \end{matrix},$$

式中: $A_{11}^{(3)}$ 为上三角矩阵.

重复上述过程, 假定已经进行了 $k-1$ 步, 得到了 Householder 变换 $H_1, H_2, \cdots, H_{k-1}$, 使得

$$A_k = H_{k-1} \cdots H_2 H_1 A_1 = \begin{bmatrix} A_{11}^{(k)} & A_{12}^{(k)} \\ O & A_{22}^{(k)} \end{bmatrix} \begin{matrix} k-1 \\ m-k+1 \\ k-1 & n-k+1 \end{matrix},$$

式中: $A_{11}^{(k)}$ 为上三角矩阵.

假定

$$A_{22}^{(k)} = [a_k, a_{k+1}, \cdots, a_n].$$

第 k 步是先对非零向量 a_k , 构造

$$\widetilde{H}_k = I_{m-k+1} - \beta_k v_k v_k^T \in \mathbb{R}^{(m-k+1) \times (m-k+1)},$$

使得 $\widetilde{H}_k a_k = \alpha_k e_1$, 这里 $\alpha_k = \|a_k\|_2$, $e_1 \in \mathbb{R}^{m-k+1}$. 令

$$H_k = \text{diag}(I_{k-1}, \widetilde{H}_k),$$

则

$$\begin{aligned} A_{k+1} &= H_k A_k = \begin{bmatrix} A_{11}^{(k)} & A_{12}^{(k)} \\ O & \widetilde{H}_k A_{22}^{(k)} \end{bmatrix} \\ &= \begin{bmatrix} A_{11}^{(k+1)} & A_{12}^{(k+1)} \\ O & A_{22}^{(k+1)} \end{bmatrix} \begin{matrix} k \\ m-k \\ k & n-k \end{matrix}, \end{aligned}$$

式中: $A_{11}^{(k+1)}$ 为上三角矩阵.

这样, 从 $k=1$ 出发, 对 A 依次进行 n 次 Householder 变换, 就可将 A 约化为上三角矩阵 \widetilde{R} , 即

$$\widetilde{R} = A_{n+1} = H_n A_n = \begin{bmatrix} A_{11}^{(n+1)} \\ O \end{bmatrix} \begin{matrix} n \\ m-n \end{matrix},$$

式中: $A_{11}^{(n+1)}$ 为上三角矩阵.

现记

$$R = A_{11}^{(n+1)}, \quad Q^T = H_n H_{n-1} \cdots H_1,$$

则

$$A = Q \begin{bmatrix} R \\ O \end{bmatrix},$$

式中: $R \in \mathbb{R}^{n \times n}$ 为上三角矩阵.

若令 $Q = [Q_1, Q_2]$, 这里 $Q_1 \in \mathbb{R}^{m \times n}$, 则有

$$A = Q_1 R.$$

证毕. □

注 2.7 若 $A \in \mathbb{R}^{n \times n}$ 为非奇异矩阵, 则由定理 5.4, A 有 QR 分解

$$A = QR,$$

式中: $Q \in \mathbb{R}^{n \times n}$ 为正交矩阵; $R \in \mathbb{R}^{n \times n}$ 为非奇异的上三角矩阵.

下面考虑计算 A 的 QR 分解的存储问题. 一般来说, 在完成 QR 分解之后 A 就不再需要, 可用它来存放 Q 和 R . 通常不必将 Q 显式地算出, 而只存放构成它的 n 个 Householder 矩阵 $H_k (k = 1, 2, \dots, n)$, 而对每个 H_k , 只需保存 v_k 和 β_k 即可. 注意到 v_k 具有如下形式:

$$v_k = (1, v_{k+1}^{(k)}, \dots, v_n^{(k)})^T \in \mathbb{R}^{m-k+1},$$

正好可以将 v_k 的第 2 到 $m - k + 1$ 个分量, 即 $v_k(2:m - k + 1)$ 存放在 A 的第 k 列对角元以下的位置上, 而 A 的上三角部分用来存放 R 的上三角部分.

综合上述讨论, 可得如下算法.

算法 2.5 (Householder 变换 QR 分解)

for $k = 1 : n$

if $k < m$

$[v, \beta] = \text{r_house}(A(k:m, k));$

$A(k:m, k:n) = (I_{m-k+1} - \beta v v^T) A(k:m, k:n);$

$d(k) = \beta;$

$A(k+1:m, k) = v(2:m - k + 1);$

end

end

注 2.8 算法 2.5 是指对于给定的矩阵 $A \in \mathbb{R}^{m \times n} (m \geq n)$, 计算 Householder 矩阵 H_1, H_2, \dots, H_n 满足: 如果 $Q = H_1 H_2 \cdots H_n$, 则 $Q^T A = R$ 是上三角形矩阵, 且 A 的上三角部分被 R 的上三角部分所覆盖, 第 k 个 Householder 向量的第 $k+1$ 到第 m 个分量存放于 $A(k+1:m, k)$, $k < m$ 的位置. 容易计算出, 该算法的运算量为 $2n^2(m - n/3)$ 个 flop.

根据上述算法可编制 MATLAB 程序如下:

```
function [A,d]=house_qr(A)
%Householder变换QR分解, 不显式计算和存储正交矩阵Q
[m,n]=size(A);
for k=1:n
    if k<m
        [v,beta]=r_house(A(k:m,k));
        I=eye(m-k+1);
        A(k:m,k:n)=(I-beta*v*v')*A(k:m,k:n);
        d(k)=beta;
        A(k+1:m,k)=v(2:m-k+1);
    end
end
```

例 2.3 利用算法 2.5 对矩阵 A 进行 QR 分解, 其中

$$A = \begin{bmatrix} 76 & 96 & 85 & 35 \\ 26 & 55 & 26 & 20 \\ 51 & 14 & 82 & 26 \\ 70 & 15 & 25 & 62 \\ 90 & 26 & 93 & 48 \end{bmatrix}.$$

解 由于算法 2.5 中没有显式计算和存储正交矩阵 Q , 故需对程序 house_qr.m 稍作修改:

```
function [Q,R]=Qhouse_qr(A)
%Householder QR分解, 显式计算并存储正交矩阵Q
[m,n]=size(A); Q=eye(m);
for k=1:n
    if k<m
        [v,beta]=r_house(A(k:m,k));
        H=eye(m-k+1)-beta*v*v';
        A(k:m,k:n)=H*A(k:m,k:n);
        Q=Q*blkdiag(eye(k-1),H);
    end
end
R=triu(A);
```

然后再在命令窗口依次运行如下代码:

```
>> A=[76,96,85,35;26,55,26,20; ...
>>    51,14,82,26;70,15,25,62;90,26,93,48];
>> [Q,R]=Qhouse_qr(A)
```

即得矩阵 A 的 QR 分解, 且满足 $\|A - QR\|_2 = 4.3169 \times 10^{-14}$.

2.2.2 Givens 变换 QR 分解

Givens 变换也可以实现矩阵 A 的 QR 分解. 先考虑 A 为非奇异实方阵的情形, 有下面的定理.

定理 2.7 设 $A \in \mathbb{R}^{n \times n}$ 为非奇异矩阵, 则存在有限个 Givens 矩阵的乘积 Q , 使得 $Q^T A$ 为可逆上三角矩阵 R , 即 $A = QR$.

证明 因为 $A = (a_{ij}) \in \mathbb{R}^{n \times n}$ 可逆, 故 A 的第 1 列 $a^{(0)} \neq 0$. 于是可以构造有限个 Givens 矩阵的乘积 $G_0 \in \mathbb{R}^{n \times n}$, 使得 $G_0 a^{(0)} = \|a^{(0)}\|_2 e_1$, 这里 $e_1 \in \mathbb{R}^n$. 记 $a_{11}^{(1)} = \|a^{(0)}\|_2$, 则有

$$G_0 A = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \cdots & a_{1n}^{(1)} \\ 0 & & & \\ \vdots & & A^{(1)} & \\ 0 & & & \end{bmatrix}.$$

易见 $A^{(1)}$ 可逆. 于是 $A^{(1)}$ 的第 1 列 $a^{(1)} \neq 0$. 故存在有限个 Givens 矩阵的乘积 $G_1 \in \mathbb{R}^{(n-1) \times (n-1)}$, 使得 $G_1 a^{(1)} = \|a^{(1)}\|_2 e_1$, 这里 $e_1 \in \mathbb{R}^{n-1}$. 记 $a_{22}^{(2)} = \|a^{(1)}\|_2$, 则有

$$G_1 A^{(1)} = \begin{bmatrix} a_{22}^{(2)} & a_{23}^{(2)} & \cdots & a_{2n}^{(2)} \\ 0 & & & \\ \vdots & & A^{(2)} & \\ 0 & & & \end{bmatrix}.$$

依此类推, 最后得到 $A^{(n-2)} \in \mathbb{R}^{2 \times 2}$ 可逆, $A^{(n-2)}$ 的第 1 列 $a^{(n-2)} \neq 0$. 那么可构造 Givens 矩阵 $G_{n-2} \in \mathbb{R}^{2 \times 2}$, 使得 $G_{n-2} a^{(n-2)} = \|a^{(n-2)}\|_2 e_1$, 这里 $e_1 \in \mathbb{R}^2$. 记 $a_{n-1,n-1}^{(n-1)} = \|a^{(n-2)}\|_2$, 则有

$$G_{n-2} A^{(n-2)} = \begin{bmatrix} a_{n-1,n-1}^{(n-1)} & a_{n-1,n}^{(n-1)} \\ 0 & a_{n,n}^{(n-1)} \end{bmatrix}.$$

令

$$Q^T = \begin{bmatrix} I_{n-2} & \\ & G_{n-2} \end{bmatrix} \cdots \begin{bmatrix} I_2 & \\ & G_2 \end{bmatrix} \begin{bmatrix} 1 & \\ & G_1 \end{bmatrix} G_0,$$

则 Q 仍为有限个 Givens 矩阵之积, 且有

$$Q^T A = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \cdots & a_{1,n-1}^{(1)} & a_{1n}^{(1)} \\ & a_{22}^{(2)} & \cdots & a_{2,n-1}^{(2)} & a_{2n}^{(2)} \\ & & \ddots & \vdots & \vdots \\ & & & a_{n-1,n-1}^{(n-1)} & a_{n-1,n}^{(n-1)} \\ & & & & a_{nn}^{(n-1)} \end{bmatrix} := R,$$

即 $A = QR$. 证毕. \square

定理 2.7 的结论可以推广到 A 是长方形的情形.

定理 2.8 设 $A \in \mathbb{R}^{m \times n} (m \geq n)$ 是列满秩的, 则有 m 阶正交矩阵 Q 及 n 阶上三角矩阵 R , 使得

$$A = Q \begin{bmatrix} R \\ O \end{bmatrix} = Q_1 R,$$

式中: Q_1 为由 Q 的前 n 列构成的矩阵.

证明 已知 A 是列满秩的, 即它的 n 个列向量线性无关, 使用与证明定理 2.7 相同的方法, 可找到有限个 m 阶 Givens 矩阵的乘积 G , 使得 $GA = \begin{bmatrix} R \\ O \end{bmatrix}$, 再令 $Q = G^T$ 即得所求. 证毕. \square

下面考虑 Givens 变换 QR 分解的实现过程. 可以用一个 5×3 阶矩阵的例子来表明其一般思想.

$$\begin{aligned} & \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix} \xrightarrow{(4,5)} \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \\ 0 & \times & \times \end{bmatrix} \xrightarrow{(3,4)} \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \\ 0 & \times & \times \\ 0 & \times & \times \end{bmatrix} \xrightarrow{(2,3)} \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ 0 & \times & \times \\ 0 & \times & \times \\ 0 & \times & \times \end{bmatrix} \xrightarrow{(1,2)} \\ & \begin{bmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & \times & \times \\ 0 & \times & \times \\ 0 & \times & \times \end{bmatrix} \xrightarrow{(4,5)} \begin{bmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & \times & \times \\ 0 & \times & \times \\ 0 & 0 & \times \end{bmatrix} \xrightarrow{(3,4)} \begin{bmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & \times & \times \\ 0 & 0 & \times \\ 0 & 0 & \times \end{bmatrix} \xrightarrow{(2,3)} \begin{bmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & 0 & \times \\ 0 & 0 & \times \\ 0 & 0 & \times \end{bmatrix} \xrightarrow{(4,5)} \\ & \begin{bmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & 0 & \times \\ 0 & 0 & \times \\ 0 & 0 & 0 \end{bmatrix} \xrightarrow{(3,4)} \begin{bmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & 0 & \times \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} := \begin{bmatrix} R_1 \\ O \end{bmatrix} := R, \end{aligned}$$

此处已标记定义所对应 Givens 变换的 2 维向量. 显然, 若 G_i 表示在约化过程中的第 i 次 Givens 变换, 则 $Q^T A = R$ 是上三角形矩阵, 其中 $Q = G_1 G_2 \cdots G_t$, 且 t 是 Givens 变换的总次数. 对于一般的 $A \in \mathbb{R}^{m \times n} (m \geq n)$, 有如下算法.

算法 2.6 (Givens 变换 QR 分解)

$[m, n] = \text{size}(A);$

for $k = 1 : n$

for $i = m : -1 : k + 1$

```

[c,s]=r_givens(A(i-1,k),A(i,k));
A(i-1:i,k:n)= $\begin{bmatrix} c & s \\ -s & c \end{bmatrix}$ A(i-1:i,k:n);
end
end

```

注 2.9 算法 2.6 是指对于给定的矩阵 $A \in \mathbb{R}^{m \times n} (m \geq n)$, 计算 Givens 矩阵 $G_1, G_2, \dots, G_t (t = \frac{1}{2}n(2m-n-1))$ 满足: 如果 $Q = G_1 G_2 \dots G_t$, 则 $Q^T A = R$ 是上三角形矩阵, 且 A 被 R 所覆盖. 容易计算出, 该算法的运算量为 $3n^2(m-n/3)$ 个 flop.

根据上述算法可编制 MATLAB 程序如下:

```

function [A]=givens_qr(A)
%Givens变换QR分解
[m,n]=size(A);
for k=1:n
    for i=m:-1:k+1
        [c,s]=r_givens(A(i-1,k), A(i,k));
        A(i-1:i, k:n)=[c, s; -s, c]*A(i-1:i, k:n);
    end
end
end

```

例 2.4 利用算法 2.6 对例 2.3 中的矩阵 A 进行 QR 分解.

解 由于算法 2.6 中没有显式计算和存储正交矩阵 Q , 故不能直接调用程序 givens_qr.m, 需对其稍作修改:

```

function [Q,R]=Qgivens_qr(A)
%Givens变换QR分解, 显式计算并存储正交矩阵Q
[m,n]=size(A); Q=eye(m);
for k=1:n
    for i=m:-1:k+1
        [c,s]=r_givens(A(i-1,k),A(i,k));
        A(i-1:i,k:n)=[c,s;-s,c]*A(i-1:i,k:n);
        Q(:,i-1:i)=Q(:,i-1:i)*[c,s;-s,c]';
    end
end
end
R=triu(A);

```

然后再在命令窗口依次运行如下代码:

```

>> A=[76,96,85,35;26,55,26,20; ...
>> 51,14,82,26;70,15,25,62;90,26,93,48];
>> [Q,R]=Qgivens_qr(A)

```

即得矩阵 A 的 QR 分解, 且满足 $\|A - QR\|_2 = 7.6458 \times 10^{-14}$.

作为 Givens 变换 QR 分解的一个应用, 下面讨论上 Hessenberg 矩阵的 QR 分解问题. 从前述讨论可知, Givens 变换 QR 分解的运算量约为 Householder 变换的 1.5 倍. 尽管如此, 用 Givens 变换对上 Hessenberg 矩阵进行 QR 分解却具有独特的优势. 上 Hessenberg 矩阵是一类具有重要应用的矩阵类, 例如第 4 章即将介绍的广义极小残量法 (GMRES) 需要求解一个上 Hessenberg 矩阵的最小二乘问题. 下面介绍上 Hessenberg 矩阵的 Givens 变换 QR 分解的具体实现.

定义 2.4 如果 n 阶矩阵的下次对角线以下的元素都为零, 则称该矩阵为上 Hessenberg 矩阵. 上 Hessenberg 矩阵的转置称为下 Hessenberg 矩阵.

根据上 Hessenberg 矩阵的特殊结构 (矩阵的每一列对角线以下的元素中只有第 1 个元素可能非零), 因而只需对每一列作一次 Givens 变换即可, 具体来说有下面的算法.

算法 2.7 (上 Hessenberg 矩阵 QR 分解)

for $k = 1 : n - 1$

$[c_k, s_k] = \text{r.givens}(A(k, k), A(k + 1, k));$

$$A(k : k + 1, k : n) = \begin{bmatrix} c_k & s_k \\ -s_k & c_k \end{bmatrix} A(k : k + 1, k : n);$$

end

注 2.10 算法 2.7 是指对于给定的上 Hessenberg 矩阵 $A \in \mathbb{R}^{n \times n}$, 计算 Givens 矩阵 G_1, G_2, \dots, G_{n-1} 满足: 如果 $Q = G_1 G_2 \cdots G_{n-1}$, 则 $Q^T A = R$ 是上三角形矩阵, 且 A 被 R 所覆盖, 其中 G_k 形如 $G_{k, k+1}(c_k, s_k)$. 容易计算出, 该算法的运算量只有 $3n^2$ 个 flop.

根据上述算法可编制 MATLAB 程序如下:

```
function [A]=hessenberg_qr(A)
%上Hessenberg矩阵的QR分解
[n]=size(A,2);
for k=1:n-1
    [c,s]=r_givens(A(k,k),A(k+1,k));
    A(k:k+1,k:n)=[c,s;-s,c]*A(k:k+1,k:n);
end
```

作为 Givens 变换 (或 Householder 变换) 的另一个应用, 下面介绍 n 阶实矩阵正交相似于上 Hessenberg 矩阵的计算问题.

定理 2.9 设 $A \in \mathbb{R}^{n \times n}$, 则存在有限个 Givens 矩阵 (或 Householder 矩阵) 的乘积 Q , 使得 $Q A Q^T$ 为上 Hessenberg 矩阵.

证明 仅讨论使用 Givens 矩阵的情形.

第 1 步, 设 $A = (a_{ij}) \in \mathbb{R}^{n \times n}$, 记 $\alpha^{(0)} = (a_{21}, \dots, a_{n1})^T \in \mathbb{R}^{n-1}$, 当 $\alpha^{(0)} = 0$ 时转入第 2 步; 当 $\alpha^{(0)} \neq 0$ 时, 构造有限个 Givens 矩阵的乘积 G_0 , 使得

$$G_0 \alpha^{(0)} = \|\alpha^{(0)}\|_2 e_1, \quad (e_1 \in \mathbb{R}^{n-1}).$$

记 $a_{21}^{(1)} = \|\alpha^{(0)}\|_2$, 则有

$$\begin{bmatrix} 1 & & \\ & G_0 & \end{bmatrix} A \begin{bmatrix} 1 & & \\ & G_0 & \end{bmatrix}^T = \begin{bmatrix} a_{11} & a_{12}^{(1)} & a_{13}^{(1)} & \cdots & a_{1n}^{(1)} \\ a_{21}^{(1)} & & & & \\ 0 & & & & \\ \vdots & & & A^{(1)} & \\ 0 & & & & \end{bmatrix}.$$

第 2 步, 设 $A^{(1)} \in \mathbb{R}^{(n-1) \times (n-1)}$, 记 $\alpha^{(1)} = (a_{32}^{(1)}, \dots, a_{n2}^{(1)})^T \in \mathbb{R}^{n-2}$, 当 $\alpha^{(1)} = 0$ 时转入第 3 步; 当 $\alpha^{(1)} \neq 0$ 时, 构造有限个 Givens 矩阵的乘积 G_1 , 使得

$$G_1 \alpha^{(1)} = \|\alpha^{(1)}\|_2 e_1, \quad (e_1 \in \mathbb{R}^{n-2}).$$

记 $a_{32}^{(2)} = \|\alpha^{(1)}\|_2$, 则有

$$\begin{bmatrix} 1 & & \\ & G_1 & \end{bmatrix} A^{(1)} \begin{bmatrix} 1 & & \\ & G_1 & \end{bmatrix}^T = \begin{bmatrix} a_{22}^{(1)} & a_{23}^{(2)} & a_{24}^{(2)} & \cdots & a_{2n}^{(2)} \\ a_{32}^{(2)} & & & & \\ 0 & & & & \\ \vdots & & & A^{(2)} & \\ 0 & & & & \end{bmatrix}.$$

继续上述过程, 直到第 $n-2$ 步:

$$A^{(n-3)} = \begin{bmatrix} a_{n-2,n-2}^{(n-3)} & a_{n-2,n-1}^{(n-3)} & a_{n-2,n}^{(n-3)} \\ a_{n-1,n-2}^{(n-3)} & a_{n-1,n-1}^{(n-3)} & a_{n-1,n}^{(n-3)} \\ a_{n,n-2}^{(n-3)} & a_{n,n-1}^{(n-3)} & a_{n,n}^{(n-3)} \end{bmatrix} \in \mathbb{R}^{3 \times 3},$$

记 $\alpha^{(n-3)} = \begin{bmatrix} a_{n-1,n-2}^{(n-3)} \\ a_{n,n-2}^{(n-3)} \end{bmatrix}$. 当 $\alpha^{(n-3)} = 0$ 时结束; 当 $\alpha^{(n-3)} \neq 0$ 时, 构造 Givens 矩阵 G_{n-3} , 使得

$$G_{n-3} \alpha^{(n-3)} = \|\alpha^{(n-3)}\|_2 e_1, \quad (e_1 \in \mathbb{R}^2).$$

记 $a_{n-1,n-2}^{(n-2)} = \|\alpha^{(n-3)}\|_2$, 则有

$$\begin{bmatrix} 1 & & \\ & G_{n-3} & \end{bmatrix} A^{(n-3)} \begin{bmatrix} 1 & & \\ & G_{n-3} & \end{bmatrix}^T = \begin{bmatrix} a_{n-2,n-2}^{(n-3)} & a_{n-2,n-1}^{(n-3)} & a_{n-2,n}^{(n-3)} \\ a_{n-1,n-2}^{(n-2)} & a_{n-1,n-1}^{(n-2)} & a_{n-1,n}^{(n-2)} \\ 0 & a_{n,n-1}^{(n-2)} & a_{n,n}^{(n-2)} \end{bmatrix}.$$

最后, 构造正交矩阵

$$Q = \begin{bmatrix} I_{n-2} & & \\ & G_{n-3} & \end{bmatrix} \cdots \begin{bmatrix} I_2 & \\ & G_1 \end{bmatrix} \begin{bmatrix} 1 & & \\ & G_0 & \end{bmatrix},$$

可使 $Q A Q^T$ 为上 Hessenberg 矩阵. 证毕. \square

推论 2.1 设 A 是 n 阶实对称矩阵, 则存在有限个 Givens 矩阵 (或 Householder 矩阵) 的乘积 Q , 使得 $Q A Q^T$ 为实对称三对角矩阵.

证明 对 A 应用定理 2.9, 则存在正交矩阵 Q , 使得 $Q A Q^T$ 为上 Hessenberg 矩阵. 由于 $A^T = A$, 又有 $Q A Q^T = (Q A Q^T)^T$ 为下 Hessenberg 矩阵, 故 $Q A Q^T$ 是三对角矩阵, 从而 $Q A Q^T$ 是实对称的三对角矩阵. 证毕. \square

2.3 线性无关向量组的正交化

本节讨论线性无关向量组的正交化. 对以这些向量为列的矩阵施以 Householder 变换或 Givens 变换作正交化, 就可以得到该线性无关向量组的正交化. 首先介绍经典的 Gram-Schmidt 正交化过程.

2.3.1 Gram-Schmidt 正交化

考虑 m 个线性无关的 n 维向量组 $\{x_i\}_{i=1}^m$ ($m \leq n$). 记列满秩 $n \times m$ 阶矩阵 $X = [x_1, x_2, \dots, x_m]$. 对向量组 $\{x_i\}_{i=1}^m$ 用 Gram-Schmidt 正交化产生相互正交的单位向量组 $\{q_i\}_{i=1}^m$, 使得

$$\text{span}\{x_1, x_2, \dots, x_k\} = \text{span}\{q_1, q_2, \dots, q_k\}, \quad k = 1, 2, \dots, m. \quad (2.10)$$

按如下方式构造 $\{q_i\}_{i=1}^m$. 令 $q_1 = x_1 / \|x_1\|_2$. 设已经构造好相互正交的单位向量组 $\{q_i\}_{i=1}^{k-1}$ 使得式 (2.10) 对 $k-1$ 成立. 利用正交性, 得

$$\tilde{q}_k = x_k - \sum_{i=1}^{k-1} r_{ik} q_i, \quad r_{is} = (x_k, q_i), \quad i = 1, 2, \dots, k-1; \quad q_k = \frac{\tilde{q}_k}{\|\tilde{q}_k\|_2}. \quad (2.11)$$

从而 q_k 与 q_1, q_2, \dots, q_{k-1} 都正交且为单位向量. 在 q_1, q_2, \dots, q_{k-1} 都已经单位正交化的前提下, 式 (2.11) 表明 x_k 去掉它在 q_i 上的分量 (投影) $r_{ik} q_i$ ($i = 1, 2, \dots, k-1$) 后剩下的部分与 q_1, q_2, \dots, q_{k-1} 都正交. 据此, 可以写出 Gram-Schmidt 正交化的算法如下:

算法 2.8 (标准 Gram-Schmidt 正交化) 给定 $n \times m$ ($m \leq n$) 阶列满秩矩阵 $X = [x_1, x_2, \dots, x_m]$, 本算法产生 $n \times m$ 阶矩阵 $Q = [q_1, q_2, \dots, q_m]$ (其列是单位正交的) 和 $m \times m$ 阶非奇异上三角矩阵 $R = (r_{ik})$.

```

 $r_{11} = \|x_1\|_2; q_1 = x_1/r_{11};$ 
for  $k = 2 : m$ 
    for  $(i = 1 : k - 1), r_{ik} = (x_k, q_i);$  end
     $\tilde{q} = x_k - \sum_{i=1}^{k-1} r_{ik} q_i;$ 
     $r_{kk} = \|\tilde{q}\|_2; q_k = \tilde{q}/r_{kk};$ 
end

```

由算法 2.8 可以得到关系式

$$x_k = \sum_{i=1}^k r_{ik} q_i, \quad k = 1, 2, \dots, m. \quad (2.12)$$

写成矩阵形式即为

$$X = QR, \quad (2.13)$$

式中: $Q = [q_1, q_2, \dots, q_m] \in \mathbb{R}^{n \times m}$, 其列是相互正交的单位向量; R 为 $m \times m$ 阶上三角矩阵, x_1, x_2, \dots, x_m 的线性无关性保证了所有的 r_{ii} ($i = 1, 2, \dots, m$) 都不为零, 因而 R 是非奇异的.

算法 2.8 表明, 通过对矩阵 X 列向量的 Gram-Schmidt 正交化, 实现了矩阵 X 的 QR 分解.

标准 Gram-Schmidt 正交化过程的 MATLAB 程序如下:

```

%标准Gram-Schmidt正交化程序-G_Schmidt1.m
function [Q,R]=G_Schmidt1(X)
[m]=size(X,2); %矩阵的列
R=zeros(m); R(1,1)=norm(X(:,1));
Q(:,1)=X(:,1)/R(1,1); %单位化
for k=2:m
    for (i=1:k-1), R(i,k)=Q(:,i)'*X(:,k); end
    qt=X(:,k);
    for (i=1:k-1), qt=qt-R(i,k)*Q(:,i); end
    R(k,k)=norm(qt); Q(:,k)=qt/R(k,k); %单位化
end

```

从算法 2.8 可以看出, 当 x_k 与前面的 x_i 接近线性相关时, r_{kk} 将接近于 0, 用它作分母会导致很大的舍入误差, q_i 之间很快失去了正交性. 一种改进就是修正的 Gram-Schmidt 正交化. 它把算法 2.8 的第 3 行和第 4 行修改为下面的循环:

$$\tilde{q} = x_k, \quad r_{ik} = (\tilde{q}, q_i), \quad \tilde{q} := \tilde{q} - r_{ik} q_i, \quad i = 1, 2, \dots, k-1. \quad (2.14)$$

在不考虑舍入误差时, 上述计算结果 \tilde{q} 与算法 2.8 中第 4 行的 \tilde{q} 是相同的. 实际上, 修正的 Gram-Schmidt 正交化利用了最新的 \tilde{q} , 有助于减少舍入误差影响.

为了保证正交性, 更可靠的方法是进行重正交化. 设已用式 (2.14) 计算了 \tilde{q} . 若 $\|\tilde{q}\|_2$ 很小 (与 $\|x_k\|_2$ 相比), 用它作分母计算 q_k 会有很大的舍入误差. 此时, 对已计算的 \tilde{q} 再正交化, 得

$$\tilde{r}_{ik} = (\tilde{q}, q_i), \quad \tilde{q} := \tilde{q} - \tilde{r}_{ik} q_i, \quad i = 1, 2, \dots, k-1. \quad (2.15)$$

由于 \tilde{q} 已与 q_i 接近正交, 故 \tilde{r}_{ik} 接近 0. 把它累加到 r_{ik} 上, 得

$$r_{ik} := r_{ik} + \tilde{r}_{ik}, \quad i = 1, 2, \dots, k-1. \quad (2.16)$$

可以写出如下的修正 Gram-Schmidt 正交化算法.

算法 2.9 (修正 Gram-Schmidt 正交化) 给定 $n \times m$ ($m \leq n$) 阶列满秩矩阵 $X = [x_1, x_2, \dots, x_m]$, 本算法产生 $n \times m$ 阶矩阵 $Q = [q_1, q_2, \dots, q_m]$ (其列是单位正交的) 和 $m \times m$ 阶非奇异上三角矩阵 $R = (r_{ik})$.

```

 $r_{11} = \|x_1\|_2; q_1 = x_1/r_{11};$ 
for  $k = 2 : m$ 
     $\tilde{q} = x_k;$ 
    for  $(i = 1 : k-1)$ 
         $r_{ik} = (\tilde{q}, q_i); \tilde{q} = \tilde{q} - r_{ik} q_i;$ 
    end
    for  $(i = 1 : k-1)$  %重正交化
         $\tilde{r}_{ik} = (\tilde{q}, q_i); \tilde{q} = \tilde{q} - \tilde{r}_{ik} q_i;$ 
         $r_{ik} = r_{ik} + \tilde{r}_{ik};$ 
    end
     $r_{kk} = \|\tilde{q}\|_2; q_k = \tilde{q}/r_{kk};$ 
end

```

再给出修正 Gram-Schmidt 正交化过程的 MATLAB 程序如下:

```

%修正Gram-Schmidt正交化程序-G_Schmidt2.m
function [Q,R]=G_Schmidt2(X)
[n,m]=size(X); %矩阵的列
R=zeros(m); Q=zeros(n,m);
R(1,1)=norm(X(:,1));
Q(:,1)=X(:,1)/R(1,1); %单位化
for k=2:m
    qt=X(:,k);
    for i=1:k-1

```



```

        R(i,k)=qt'*Q(:,i);
        qt=qt-R(i,k)*Q(:,i); %对剩余向量进行修正
    end
    for i=1:k-1 %重正交化
        rt=qt'*Q(:,i);
        qt=qt-rt*Q(:,i);
        R(i,k)=R(i,k)+rt;
    end
    R(k,k)=norm(qt);
    Q(:,k)=qt/R(k,k); %对本次得到的向量单位化
end

```

例 2.5 用 Gram-Schmidt 方法对矩阵 A 的列向量进行正交化, 其中

$$A = \begin{bmatrix} 9 & 1 & 4 & 5 \\ 6 & 2 & 1 & 3 \\ 4 & 1 & 9 & 9 \\ 5 & 2 & 9 & 4 \\ 4 & 2 & 5 & 1 \end{bmatrix}.$$

解 分别用标准 Gram-Schmidt 方法和修正 Gram-Schmidt 方法对矩阵 A 列向量进行正交化. 在命令窗口输入:

```

>> A=[9,1,4,5;6,2,1,3;4,1,9,9;5,2,9,4;4,2,5,1];
>> [Q1,R1]=G_Schmidt1(A);err1=norm(A-Q1*R1)
>> [Q2,R2]=G_Schmidt2(A);err2=norm(A-Q2*R2)

```

运行后即得结果, 且分别满足 $\|A - Q_1 R_1\|_2 = 1.0175 \times 10^{-15}$ 和 $\|A - Q_2 R_2\|_2 = 2.2204 \times 10^{-16}$.

2.3.2 Householder 正交化

2.2 节中介绍了 Householder 变换可以对矩阵 $A \in \mathbb{R}^{m \times n} (m \geq n)$ 进行 QR 分解. 这一过程也可以对任意 m 个线性无关的 n 维向量组进行. 与 Gram-Schmidt 正交化相比, 它具有更好的数值稳定性. 给定列满秩矩阵 $X = [x_1, x_2, \dots, x_m]$, 这里假设 $m \leq n$. 令 $X_1 = X$, 构造 n 阶 Householder 矩阵 H_1 使得 $X_2 = H_1 X_1$ 的第 1 列的后 $n-1$ 个分量全部为 0. 一般地, 假设构造了 $k-1$ 个 Householder 矩阵 $\{H_i\}_{i=1}^{k-1}$, 使得

$$X_k = H_{k-1} H_{k-2} \cdots H_1 X_1 = \begin{bmatrix} * & * & * & * & \cdots & * \\ & \ddots & * & * & \cdots & * \\ & & * & * & \cdots & * \\ & & & x_{kk} & \cdots & x_{km} \\ & & & \vdots & \ddots & \vdots \\ & & & x_{nk} & \cdots & x_{nm} \end{bmatrix}. \quad (2.17)$$

由于 x_1, \dots, x_m 线性无关, 故 X_k 的前 k 列线性无关. 故 $\tilde{x}_k = (x_{kk}, x_{k+1,k}, \dots, x_{nk})^T$ 不是零向量. 构造 $n-k+1$ 阶 Householder 矩阵 $\widetilde{H}_k = I_{n-k+1} - \beta_k v_k v_k^T$, 使得 $\widetilde{H}_k \tilde{x}_k$ 的后 $n-k$ 个分量为 0. 把 n 阶 Householder 矩阵 $H_k = \text{diag}(I_{k-1}, \widetilde{H}_k)$ 作用到 X_k 后得到

$$X_{k+1} = H_k X_k = H_k H_{k-1} \cdots H_1 X_1 = \begin{bmatrix} * & * & * & * & * & \cdots & * \\ & \ddots & * & * & * & \cdots & * \\ & & * & * & * & \cdots & * \\ & & & \otimes & \otimes & \cdots & \otimes \\ & & & 0 & \otimes & \cdots & \otimes \\ & & & \vdots & \vdots & \ddots & \vdots \\ & & & 0 & \otimes & \cdots & \otimes \end{bmatrix}. \quad (2.18)$$

X_{k+1} 与 X_k 相比, 改变的仅仅是 X_k 右下角的 $n-k+1$ 阶子矩阵, 并产生如式 (2.18) 右端矩阵中位置的零元素. 上述过程进行 m 次之后, 就产生了 $n \times m$ 阶矩阵

$$X_{m+1} = H_m H_{m-1} \cdots H_1 X_1 = \begin{bmatrix} R \\ O \end{bmatrix},$$

式中: R 为 m 阶上三角矩阵. 令

$$H = H_m H_{m-1} \cdots H_1, \quad Q = H^T \begin{bmatrix} I_m \\ O \end{bmatrix} = [q_1, q_2, \dots, q_m] \in \mathbb{R}^{n \times m},$$

则

$$X = H^T \begin{bmatrix} R \\ O \end{bmatrix} = H^T \begin{bmatrix} I_m \\ O \end{bmatrix} R = QR,$$

这就对一般矩阵实现了 QR 分解.

下面考虑这一分解算法的程序实现. 注意到

$$H_k e_i^{(n)} = e_i^{(n)}, \quad k > i,$$

式中: $e_i^{(n)}$ 为 n 维坐标向量, 故 $Q = [q_1, q_2, \dots, q_m]$ 满足

$$q_i = Q e_i^{(n)} = H_1 H_2 \cdots H_m e_i^{(n)} = H_1 H_2 \cdots H_i e_i^{(n)}.$$

为了描述整个算法, 再引入 X_{m+1} 的第 k 列为

$$\begin{aligned} r_k &= H_m H_{m-1} \cdots H_1 X_1 e_k = H_m H_{m-1} \cdots H_1 x_k \\ &= H_k H_{k-1} \cdots H_1 x_k, \quad 1 \leq k \leq m, \end{aligned} \quad (2.19)$$

它就是 X_{k+1} 的第 k 列. 显然, 在计算出 r_k 之后, 在以后的各次正交化过程中, 它不再改变. 下面的算法给出了这一正交化过程的详细步骤.

算法 2.10 (Householder 正交化) 给定 $n \times m$ 阶列满秩矩阵 $X = [x_1, x_2, \dots, x_m]$, 本算法产生 $n \times m$ 阶矩阵 $X_{m+1} = [r_1, r_2, \dots, r_m]$ 和 $Q = [q_1, q_2, \dots, q_m]$ (其列是单位正交的, 如果不需要 Q 则不需计算). X_{m+1} 的前 m 行构成 m 阶上三角矩阵 R .

① $r_1 = x_1$;

for $k = 1 : m$

② if $(k > 1)$, 计算 $r_k = H_{k-1}H_{k-2} \cdots H_1 x_k$; end

③ 根据 r_k 的后 $n - k + 1$ 个分量构成的向量 \tilde{x} , 计算 \widetilde{H}_k 满足

$$\widetilde{H}_k \tilde{x} = \eta_k e_1^{(n-k+1)}; \text{ 令 } H_k = \text{diag}(I_{k-1}, \widetilde{H}_k);$$

④ 计算 $r_k = H_k r_k$;

⑤ 计算 $q_k = H_1 H_2 \cdots H_k e_k$;

end

注 2.11 算法 2.10 的第 ② 步可采取如下方式实现:

$z = x_k$;

for $(i = 1 : k - 1)$, $z = H_i z$; end

$r_k = z$;

(2.20)

类似地, 第 ⑤ 步可采取如下方式实现:

$z = e_k$;

for $(i = k : -1 : 1)$, $z = H_i z$; end

$q_k = z$;

现在来看一下算法 2.10 的程序执行. 当 $k = 1$ 时, 根据 $r_1 = x_1$ 的后 n 个分量构成 \tilde{x}_1 (就是 n 维列向量 x_1), 构造 $\widetilde{H}_1 = H_1$, 计算 X_{m+1} 的第 1 列 $r_1 := H_1 x_1$, $k = 1$ 迭代结束. 当 $k = 2$ 时, 计算 $r_2 := H_1 x_2$, 根据 r_2 的后 $n - 1$ 个分量构成 \tilde{x}_2 , 构造 \widetilde{H}_2 和 $H_2 = \text{diag}(I_1, \widetilde{H}_2)$. 计算 X_{m+1} 的第 2 列 $r_2 := H_2 r_2 = H_2 H_1 x_1$. 至此, 确实 r_2 已经计算了. 注意算法 2.10 的第 ② 步计算的 r_k 是 X_k 的第 k 列, 第 ④ 步计算的 r_k 才是 X_{k+1} 的第 k 列, 亦即 X_{m+1} 的第 k 列.

此外, 该算法要保存 $\widetilde{H}_k = I_{n-k+1} - \beta_k v_k v_k^T$ 中的数 β_k 和 $n - k + 1$ 维列向量 v_k . 在式 (2.20) $H_i z$ 的计算中, 由于 $H_i = \text{diag}(I_{i-1}, \widetilde{H}_i)$, 所以 $H_i z$ 的前 $i - 1$ 个分量与 z 的前 $s - 1$ 个分量相同, 而 $H_i z$ 的后 $n - i + 1$ 个分量构成的列向量就是 \widetilde{H}_i 与 z 后 $n - i + 1$ 个分量构成的列向量的乘积.

算法 2.10 的 MATLAB 程序如下:

```
%Householder正交化程序-House_orth.m
function [Q,R]=House_orth(X)
[n,m]=size(X); %矩阵的行和列
R(:,1)=X(:,1); E=eye(n); Y=cell(m,1);
```

```

for k=1:m
    if k>1
        z=X(:,k);
        for i=1:k-1
            v=Y{i}; %取出来
            z=blkdiag(eye(i-1),eye(n-i+1)-b(i)*v*v')*z;
        end
        R(:,k)=z;
    end
    [v,beta]=r_house(R(k:n,k));
    H=blkdiag(eye(k-1),eye(n-k+1)-beta*v*v');
    R(:,k)=H*R(:,k);
    b(k)=beta; Y{k}=v; %存起来
    z=E(:,k);
    for i=k:-1:1
        v=Y{i}; %取出来
        z=blkdiag(eye(i-1),eye(n-i+1)-b(i)*v*v')*z;
    end
    Q(:,k)=z;
end
R=R(1:m,:);

```

例 2.6 用 Householder 正交化方法对例 2.5 中矩阵 A 的列向量进行正交化.

解 利用程序 House_orth.m, 在命令窗口输入:

```

A=[9,1,4,5;6,2,1,3;4,1,9,9;5,2,9,4;4,2,5,1];
[Q,R]=House_orth(A); err=norm(A-Q*R)

```

运行后即得结果, 且满足 $\|A - QR\|_2 = 4.3565 \times 10^{-15}$.

2.4 Krylov 子空间及其正交化

本节阐述与 Krylov 子空间有关的一些基本概念和重要性质, 并且给出计算其正交基的几个实用的算法.

2.4.1 Krylov 子空间

首先来介绍一下 Krylov 子空间的有关概念和性质. 考虑这样一个问题: 假如并不明确地知道矩阵 A , 而只知道对任意给定向量 x 可以产生向量 $y = Ax$ 的机制, 那么可以从中得到 A 的一些什么信息呢?

当然, 在这种情况下, 若选定一个向量 v , 便可得到一个序列

$$v_0 = v, v_1 = Av_0 = Av, v_2 = Av_1 = A^2v, \dots, v_{k+1} = Av_k = A^{k+1}v, \dots \quad (2.21)$$

假定 A 的特征多项式为

$$p_A(\lambda) = \det(\lambda I - A) = \lambda^n + \alpha_{n-1}\lambda^{n-1} + \cdots + \alpha_1\lambda + \alpha_0,$$

则由 Cayley-Hamilton 定理可知

$$v_n + \alpha_{n-1}v_{n-1} + \cdots + \alpha_1v_1 + \alpha_0v_0 = 0,$$

从而有

$$Vp = -v_n, \quad (2.22)$$

式中:

$$V = [v_0, v_1, \cdots, v_{n-1}], \quad p = (\alpha_0, \alpha_1, \cdots, \alpha_{n-1})^T.$$

若 V 非奇异, 则由方程组 (2.22) 可唯一地确定特征多项式的系数构成的向量 p , 从而可求得特征多项式, 进而得到 A 的所有特征值. 这就是说, 在所给定的条件下, 能得到 A 所有的特征值的信息.

上述求 A 的特征多项式的方法, 是在 1931 年首先由 Krylov 提出的. 因此, 后人就将序列式 (2.21) 称为 Krylov 序列, 而将子空间

$$\mathcal{K}_k(A, v) = \text{span}\{v, Av, \cdots, A^{k-1}v\} \quad (2.23)$$

称为 Krylov 子空间, 将矩阵

$$K_k(A, v) = [v, Av, \cdots, A^{k-1}v] \quad (2.24)$$

称为 Krylov 矩阵.

由 Krylov 子空间的定义, 容易验证其有如下的性质.

定理 2.10 假定 $A \in \mathbb{R}^{n \times n}$ 和 $v \in \mathbb{R}^n$ 已经给定, 其中 $v \neq 0$, 则 Krylov 子空间有如下性质:

(1) Krylov 子空间序列满足

$$\mathcal{K}_k(A, v) \subset \mathcal{K}_{k+1}(A, v), \quad A\mathcal{K}_k(A, v) \subset \mathcal{K}_{k+1}(A, v).$$

(2) 对任意的非零实数 α , 有

$$\mathcal{K}_k(A, v) = \mathcal{K}_k(\alpha A, v) = \mathcal{K}_k(A, \alpha v) \quad (\text{伸缩不变性}).$$

(3) 对任意的实数 μ , 有

$$\mathcal{K}_k(A, v) = \mathcal{K}_k(A - \mu I, v) \quad (\text{位移不变性}).$$

(4) 对任意的非奇异矩阵 $W \in \mathbb{R}^{n \times n}$, 有

$$\mathcal{K}_k(W^{-1}AW, W^{-1}v) = W^{-1}\mathcal{K}_k(A, v).$$

(5) 若记次数不超过 $k-1$ 的实系数多项式的全体为 \mathcal{P}_{k-1} , 则 $\mathcal{K}_k(A, v)$ 有如下表

示

$$\mathcal{K}_k(A, v) = \{p(A)v : p \in \mathcal{P}_{k-1}\}.$$

注 2.12 定理 2.10 的第 2 条是说, 对矩阵 A 和向量 v 乘以任意非零实数, 并不改变其生成的 Krylov 子空间, 称 Krylov 子空间的这一性质为伸缩不变性. 第 3 条是说, 对矩阵进行位移, 并不改变其生成的 Krylov 子空间, 称这一性质为位移不变性. 第 4 条性质给出了在矩阵的相似变换下相应的 Krylov 子空间之间的关系, 这对处理与 Krylov 子空间有关的理论时十分有用. 第 5 条给出了 Krylov 子空间一个重要性质, 即其每个元都可看作是某个矩阵多项式作用到初始向量 v 上而得到的.

假设 (λ, v) 是 A 的特征对, 则有 $A^i v = \lambda^i v$, 从而有

$$\mathcal{K}_k(A, v) = \mathcal{K}_1(A, v), \quad k = 1, 2, \dots$$

换句话说, 此时在 Krylov 序列中, 从第 1 项之后就没有任何新的向量产生, 即它终止于第一项. 一般来讲, Krylov 序列在第 ℓ 项终止, 是指 ℓ 为满足 $\mathcal{K}_{\ell+1}(A, v) = \mathcal{K}_\ell(A, v)$ 的最小整数.

定理 2.11 若 Krylov 序列在第 ℓ 项终止, 则 $\mathcal{K}_\ell(A, v)$ 是 A 的一个 ℓ 维不变子空间. 反过来, 若 v 属于 A 的一个 m 维不变子空间, 则必存在一个 $\ell \leq m$, 使得对应的 Krylov 序列在第 ℓ 项终止.

证明 若 Krylov 序列在第 ℓ 项终止, 则对任意 $x \in \mathcal{K}_\ell(A, v)$, 有

$$Ax \in \mathcal{K}_{\ell+1}(A, v) = \mathcal{K}_\ell(A, v),$$

即 $\mathcal{K}_\ell(A, v)$ 是 A 的不变子空间. 此外, 由于 ℓ 是使得 $\mathcal{K}_{\ell+1}(A, v) = \mathcal{K}_\ell(A, v)$ 的最小整数, 故必有 $v, Av, \dots, A^{\ell-1}v$ 这 ℓ 个向量是线性无关的, 从而就有 $\dim \mathcal{K}_\ell(A, v) = \ell$.

反过来, 若 $v \in \mathcal{X}_m$, 这里 \mathcal{X}_m 是 A 的 m 维不变子空间, 则对任意的整数 k , 必有 $A^k v \in \mathcal{X}_m$. 故必存在一个 $\ell \leq m$, 使得 $v, Av, \dots, A^{\ell-1}v$ 线性无关, 但 $v, Av, \dots, A^{\ell-1}v, A^\ell v$ 线性相关, 从而有 $\mathcal{K}_{\ell+1}(A, v) = \mathcal{K}_\ell(A, v)$. 当然这样的 ℓ 必须是使得上述等式成立的最小整数. 因此, 对应的 Krylov 序列必在第 ℓ 项终止. 证毕. \square

由定理 2.11 可知, 当 Krylov 序列在第 ℓ 项终止时, 则 $\mathcal{K}_\ell(A, v)$ 中就包含了 A 的一些特征向量. 这对完成诸多的计算任务是十分有利的, 但也有不利的一面. 例如, 如果 ℓ 很小, 那么 $\mathcal{K}_\ell(A, v)$ 中可能并不含有所需要的信息, 这就会给计算带来困难, 必须重新开始.

2.4.2 Arnoldi 正交分解

在实际应用 Krylov 子空间时, 需要它的一些基向量来描述它. 当然, 定义它的向量组应该是它的一组天然的基. 然而这组基在实际计算时是没有什么用处的, 因为随着 k 的增加, Krylov 矩阵会变得越来越病态. 因此, 为了实际计算的需要, 必须寻找其他对扰动并不敏感的基向量.

设 $K_{k+1}(A, v) = [v, Av, \dots, A^k v]$ 是列满秩的, 并假定它的 QR 分解为

$$K_{k+1}(A, v) = V_{k+1} R_{k+1}, \quad (2.25)$$

式中: $V_{k+1} = [v_1, v_2, \dots, v_{k+1}] \in \mathbb{R}^{n \times (k+1)}$ 满足 $V_{k+1}^T V_{k+1} = I_{k+1}$, R_{k+1} 为非奇异的上三角矩阵.

通常称 V_{k+1} 为 $K_{k+1}(A, v)$ 的 QR 分解的 Q 因子, 简称为 $K_{k+1}(A, v)$ 的 Q 因子. 将 V_{k+1} 和 R_{k+1} 作如下分块:

$$V_{k+1} = [V_k, v_{k+1}], \quad R_{k+1} = \begin{bmatrix} R_k & r_{k+1} \\ 0 & \rho_{k+1, k+1} \end{bmatrix}.$$

再注意到 $K_{k+1}(A, v) = [K_k(A, v), A^k v]$, 由式 (2.25), 得

$$K_k(A, v) = V_k R_k, \quad (2.26)$$

即 V_k 就是 $K_k(A, v)$ 的 QR 分解的 Q 因子. 从而有

$$\mathcal{K}_k(A, v) = \mathcal{R}(V_k) = \text{span}\{v_1, v_2, \dots, v_k\},$$

即 V_k 的列向量就构成了 Krylov 子空间 $\mathcal{K}_k(A, v)$ 的一组标准正交基.

这样自然想到利用 $K_k(A, v)$ 的 QR 分解来产生 $\mathcal{K}_k(A, v)$ 的一组正交基. 然而这一方法是不可取的, 因为 $K_k(A, v)$ 是十分病态的, 在计算机上形成 $K_k(A, v)$ 就会引起较大的误差. 因此, 需要寻求其他方法来计算 V_k .

事实上, 由式 (2.25) 和式 (2.26) 可以导出如下结果.

定理 2.12 给定矩阵 $A \in \mathbb{R}^{n \times n}$ 和向量 $v \in \mathbb{R}^n$, 并假定 $K_{k+1}(A, v)$ 是列满秩的, 而且其 QR 分解的 Q 因子为 V_{k+1} , 则必有一个 $(k+1) \times k$ 阶不可约上 Hessenberg 矩阵 \widetilde{H}_k , 使得

$$AV_k = V_{k+1} \widetilde{H}_k, \quad (2.27)$$

式中: V_k 为由 V_{k+1} 的前 k 列构成的矩阵. 反过来, 若有

$$V_{k+1} = [V_k, v_{k+1}] \in \mathbb{R}^{n \times (k+1)}, \quad V_{k+1}^T V_{k+1} = I_{k+1},$$

以及不可约的上 Hessenberg 矩阵 \widetilde{H}_k 使得式 (2.27) 成立, 则 V_{k+1} 必是 $K_{k+1}(A, v_1)$ 的 QR 分解的 Q 因子, 这里 $v_1 = V_{k+1} e_1$ 是 V_{k+1} 的第 1 列.

证明 设 $V_{k+1} = [V_k, v_{k+1}]$ 是 $K_{k+1}(A, v)$ 的 QR 分解的 Q 因子, 则由式 (2.26) 得

$$[Av, \dots, A^k v] = AK_k(A, v) = AV_k R_k. \quad (2.28)$$

再由式 (2.25) 得

$$[v, AK_k(A, v)] = K_{k+1}(A, v) = V_{k+1} R_{k+1}.$$

比较上式两边的后 k 列, 得

$$AK_k(A, v) = V_{k+1} \widehat{H}_k, \quad (2.29)$$

式中: \widehat{H}_k 为由 R_{k+1} 的后 k 列构成的 $(k+1) \times k$ 阶上 Hessenberg 矩阵.

注意: R_{k+1} 非奇异蕴含着 \widehat{H}_k 是不可约的. 将式 (2.28) 与式 (2.29) 相结合, 得

$$AV_k R_k = V_{k+1} \widehat{H}_k \implies AV_k = V_{k+1} \widetilde{H}_k,$$

式中: $\widetilde{H}_k = \widehat{H}_k R_k^{-1}$.

由于 \widehat{H}_k 是不可约的上 Hessenberg 矩阵, 而 R_k 是非奇异的上三角矩阵, 因此 \widetilde{H}_k 必为不可约的上 Hessenberg 矩阵.

反过来, 若有式 (2.27) 成立, 比较式 (2.27) 两边的各列, 得

$$Av_j = \sum_{i=1}^{j+1} h_{ij} v_i, \quad j = 1, 2, \dots, k.$$

于是有

$$\begin{aligned} v_2 &= \frac{1}{h_{21}}(Av_1 - h_{11}v_1) = \rho_{12}v_1 + \rho_{22}Av_1, \quad \rho_{22} = \frac{1}{h_{21}} \neq 0; \\ v_3 &= \frac{1}{h_{32}}(Av_2 - h_{12}v_1 - h_{22}v_2) \\ &= \frac{1}{h_{32}}(A - h_{22}I)(\rho_{12}v_1 + \rho_{22}Av_1) - \frac{h_{12}}{h_{32}}v_1 \\ &= \rho_{13}v_1 + \rho_{23}Av_1 + \rho_{33}A^2v_1, \quad \rho_{33} = \frac{\rho_{22}}{h_{32}} \neq 0. \end{aligned}$$

如此下去, 得

$$v_j = \rho_{1j}v_1 + \rho_{2j}Av_1 + \dots + \rho_{jj}A^{j-1}v_1, \quad \rho_{jj} = \frac{\rho_{j-1,j-1}}{h_{j,j-1}} \neq 0, \quad 2 \leq j \leq k+1. \quad (2.30)$$

令

$$R^{-1} = \begin{bmatrix} 1 & \rho_{12} & \rho_{13} & \cdots & \rho_{1,k+1} \\ & \rho_{22} & \rho_{23} & \cdots & \rho_{2,k+1} \\ & & \rho_{33} & \cdots & \rho_{3,k+1} \\ & & & \ddots & \vdots \\ & & & & \rho_{k+1,k+1} \end{bmatrix},$$

则将式 (2.30) 写成矩阵形式, 并且再补充第 1 列 v_1 , 得

$$[v_1, v_2, \dots, v_{k+1}] = [v_1, Av_1, \dots, A^k v_1] R^{-1},$$

即

$$K_{k+1}(A, v_1) = V_{k+1} R,$$

这表明 V_{k+1} 是 $K_{k+1}(A, v_1)$ 的 QR 分解的 Q 因子. 证毕. \square

基于这一结果, 称形如式 (2.27) 的公式为一个长度为 k 的 Arnoldi 分解; 若其中的 \widetilde{H}_k 是不可约的, 则称这一分解是不可约的, 否则称为可约的.

由定理 2.12 可知, 只需求得分解式 (2.27), 即 $K_k(A, v)$ 的 QR 分解的 Q 因子, 就能得到 Krylov 子空间 $K_k(A, v)$ 的一组标准正交基. 下面计算分解式 (2.27).

将 \widetilde{H}_k 分块为

$$\widetilde{H}_k = \begin{bmatrix} H_k \\ \beta_k e_k^T \end{bmatrix},$$

则 Arnoldi 分解式 (2.27) 可改写为

$$\mathbf{A}\mathbf{V}_k = \mathbf{V}_k\mathbf{H}_k + \beta_k\mathbf{v}_{k+1}\mathbf{e}_k^T. \quad (2.31)$$

这种形式有时使用起来更方便. 在式 (2.31) 两边左乘 \mathbf{V}_k^T , 并且注意到 \mathbf{V}_k 与 \mathbf{v}_{k+1} 的正交性, 即有 $\mathbf{H}_k = \mathbf{V}_k^T\mathbf{A}\mathbf{V}_k$. 通常称矩阵 \mathbf{H}_k 为 \mathbf{A} 关于 \mathbf{V}_k 的 Rayleigh 商. 事实上, 式 (2.31) 提供了一种由 \mathbf{V}_k 来计算 \mathbf{v}_{k+1} 的方法. 比较式 (2.31) 两边矩阵的最后一列, 得

$$\mathbf{A}\mathbf{v}_k = \mathbf{V}_k\mathbf{h}_k + \beta_k\mathbf{v}_{k+1}, \quad (2.32)$$

式中: \mathbf{h}_k 为 \mathbf{H}_k 的最后一列.

由于 \mathbf{V}_{k+1} 是列正交的, 故在式 (2.32) 两边左乘 \mathbf{V}_k^T , 得

$$\mathbf{h}_k = \mathbf{V}_k^T\mathbf{A}\mathbf{v}_k. \quad (2.33)$$

又由式 (2.32), 得

$$\beta_k\mathbf{v}_{k+1} = \mathbf{A}\mathbf{v}_k - \mathbf{V}_k\mathbf{h}_k. \quad (2.34)$$

于是, 便有

$$\beta_k = \|\mathbf{A}\mathbf{v}_k - \mathbf{V}_k\mathbf{h}_k\|_2, \quad \mathbf{v}_{k+1} = (\mathbf{A}\mathbf{v}_k - \mathbf{V}_k\mathbf{h}_k)/\beta_k. \quad (2.35)$$

这本质上是一个求向量 $\mathbf{A}\mathbf{v}_k$ 在 $\mathcal{R}(\mathbf{V}_k)^\perp$ 上的正交投影的 Gram-Schmidt 的正交化过程. 从对 Gram-Schmidt 正交化过程的研究可知, 为了保证其数值稳定性, 应该使用修正的 Gram-Schmidt 正交化过程. 记

$$\mathbf{h}_k = (h_{1k}, h_{2k}, \dots, h_{kk})^T, \quad h_{k+1,k} = \beta_k, \quad \mathbf{w} = \mathbf{A}\mathbf{v}_k, \quad \mathbf{r} = \mathbf{A}\mathbf{v}_k - \mathbf{V}_k\mathbf{h}_k.$$

利用这些记号, 将式 (2.33) 和式 (2.34) 换一种表达方式即为

$$h_{ik} = \mathbf{v}_i^T\mathbf{w}, \quad i = 1, 2, \dots, k, \quad \mathbf{r} = \mathbf{w} - \sum_{i=1}^k h_{ik}\mathbf{v}_i.$$

上式是说, 经典的 Gram-Schmidt 正交化过程, 是将 \mathbf{w} 在每个 \mathbf{v}_i 上的投影 h_{ik} 都算好之后, 再从 \mathbf{w} 中一并减去 \mathbf{w} 在 $\mathcal{R}(\mathbf{V}_k)^\perp$ 上的正交投影. 而修正的 Gram-Schmidt 正交化过程是, 当算好 \mathbf{w} 在 \mathbf{v}_1 上的正交投影 $h_{1k} = \mathbf{v}_1^T\mathbf{w}$ 后, 马上从 \mathbf{w} 中减去 $h_{1k}\mathbf{v}_1$, 即计算 $\mathbf{w}_1 = \mathbf{w} - h_{1k}\mathbf{v}_1$. 然后, 再用 $h_{2k} = \mathbf{v}_2^T\mathbf{w}_1$ 确定 h_{2k} , 而不是用 $h_{2k} = \mathbf{v}_2^T\mathbf{w}$ 确定. 当然, 从理论上讲, 二者是等价的, 即 $h_{2k} = \mathbf{v}_2^T\mathbf{w} = \mathbf{v}_2^T\mathbf{w}_1$. 然后再从 \mathbf{w}_1 中减去 $h_{2k}\mathbf{v}_2$, 即计算 $\mathbf{w}_2 = \mathbf{w}_1 - h_{2k}\mathbf{v}_2$, 依此类推. 数值实验证明, 修正的 Gram-Schmidt 正交化过程要比经典的 Gram-Schmidt 过程数值性态好得多.

综合上面的讨论, 可得如下算法.

算法 2.11 (Arnoldi 正交分解) 给定矩阵 $\mathbf{A} \in \mathbb{R}^{n \times n}$, 向量 $\mathbf{v}_1 \in \mathbb{R}^n$ ($\|\mathbf{v}_1\|_2 = 1$) 和正整数 k , 本算法计算一个长度为 k 的 Arnoldi 分解: $\mathbf{A}\mathbf{V}_k = \mathbf{V}_k\mathbf{H}_k + h_{k+1,k}\mathbf{v}_{k+1}\mathbf{e}_k^T$.

for $j = 1 : k$

$\mathbf{w} = \mathbf{A}\mathbf{v}_j$;

```

for  $i = 1 : j$ 
     $h_{ij} = \mathbf{v}_i^T \mathbf{w}; \mathbf{w} = \mathbf{w} - h_{ij} \mathbf{v}_i;$ 
end
 $h_{j+1,j} = \|\mathbf{w}\|_2;$ 
if  $h_{j+1,j} = 0$ 
    stop
else
     $\mathbf{v}_{j+1} = \mathbf{w} / h_{j+1,j};$ 
end
end

```

注 2.13 算法 2.11 只在计算 \mathbf{w} 时用到 \mathbf{A} 与一个向量作乘积, 这使得该算法可充分利用 \mathbf{A} 的稀疏性和其所具有的特殊结构. 此外, 若 Arnoldi 过程中途中断, 即计算过程中出现了 $h_{j+1,j} = 0$, 则这表明 $\mathcal{K}_j(\mathbf{A}, \mathbf{v}_1)$ 已经是 \mathbf{A} 的不变子空间. 在很多情况下, 这是有利的.

这里需要指出的是, 数值实验表明, 算法 2.11 在实际使用时, \mathbf{v}_j 之间的正交性很快就会损失掉. 解决这一问题的一种方法就是在计算过程中使用重正交化技术 (即在算法中再重复执行一次 Gram-Schmidt 正交化过程). 这就得到如下的算法.

算法 2.12 (重正交化 Arnoldi 分解) 给定矩阵 $\mathbf{A} \in \mathbb{R}^{n \times n}$, 向量 $\mathbf{v}_1 \in \mathbb{R}^n$ 满足 $\|\mathbf{v}_1\|_2 = 1$ 以及正整数 k , 本算法计算一个长度为 k 的 Arnoldi 分解: $\mathbf{A}\mathbf{V}_k = \mathbf{V}_k\mathbf{H}_k + h_{k+1,k}\mathbf{v}_{k+1}\mathbf{e}_k^T$, 并且在计算过程中使用重正交化技术.

```

for  $j = 1 : k$ 
     $\mathbf{w} = \mathbf{A}\mathbf{v}_j;$ 
    for  $i = 1 : j$ 
         $h_{ij} = \mathbf{v}_i^T \mathbf{w}; \mathbf{w} = \mathbf{w} - h_{ij} \mathbf{v}_i;$ 
    end
    for  $i = 1 : j$  (重正交化)
         $s = \mathbf{v}_i^T \mathbf{w}; h_{ij} = h_{ij} + s; \mathbf{w} = \mathbf{w} - s\mathbf{v}_i;$ 
    end
     $h_{j+1,j} = \|\mathbf{w}\|_2;$ 
    if  $h_{j+1,j} = 0$ 
        stop
    else
         $\mathbf{v}_{j+1} = \mathbf{w} / h_{j+1,j};$ 
    end
end
end

```

2.4.3 Lanczos 正交分解

当 $A \in \mathbb{R}^{n \times n}$ 是对称矩阵时, 在 Arnoldi 分解中, 其关于 V_k 的 Rayleigh 商 $H_k = V_k^T A V_k$ 就是一个对称三对角矩阵. 这样对应的 Arnoldi 分解就变成

$$A V_k = V_k T_k + \beta_k v_{k+1} e_k^T, \quad (2.36)$$

式中:

$$T_k = \begin{bmatrix} \alpha_1 & \beta_1 & & & \\ \beta_1 & \alpha_2 & & & \\ & & \ddots & & \\ & & & \ddots & \beta_{k-1} \\ & & & \beta_{k-1} & \alpha_k \end{bmatrix}.$$

此时, 称式 (2.36) 为一个长度为 k 的 Lanczos 分解.

比较式 (2.36) 两边的各列, 得

$$A v_1 = \alpha_1 v_1 + \beta_1 v_2,$$

$$A v_j = \beta_{j-1} v_{j-1} + \alpha_j v_j + \beta_j v_{j+1}, \quad j = 2, \dots, k.$$

于是有

$$\alpha_1 = v_1^T A v_1,$$

$$\beta_1 = \|A v_1 - \alpha_1 v_1\|_2, \quad v_2 = (A v_1 - \alpha_1 v_1) / \beta_1,$$

$$\alpha_j = v_j^T A v_j, \quad r_j = A v_j - \alpha_j v_j - \beta_{j-1} v_{j-1},$$

$$\beta_j = \|r_j\|_2, \quad v_{j+1} = r_j / \beta_j, \quad j = 2, \dots, k.$$

这样, 就得到了如下的算法.

算法 2.13 (Lanczos 方法) 给定矩阵 $A \in \mathbb{R}^{n \times n}$ 满足 $A^T = A$, 向量 $v_1 \in \mathbb{R}^n$ 满足 $\|v_1\|_2 = 1$ 以及正整数 k , 本算法计算一个长度为 k 的 Lanczos 分解: $A V_k = V_k T_k + \beta_k v_{k+1} e_k^T$.

$$\beta_0 = 0; \quad v_0 = 0;$$

for $j = 1 : k$

$$w = A v_j; \quad \alpha_j = v_j^T w;$$

$$w = w - \alpha_j v_j - \beta_{j-1} v_{j-1};$$

$$\beta_j = \|w\|_2;$$

if $\beta_j = 0$

stop

else

$$v_{j+1} = w / \beta_j;$$

end

end

当然在实际应用时, 该算法所产生的 \mathbf{v}_j 也将很快的失去它们之间的正交性. 弥补这一损失的方法仍然是重正交化方法. 其实重正交化就是在算法中的 $\mathbf{w} = \mathbf{w} - \alpha_j \mathbf{v}_j - \beta_{j-1} \mathbf{v}_{j-1}$ 之后, 再加入一句

$$\mathbf{w} = \mathbf{w} - \sum_{i=1}^{j-1} (\mathbf{v}_i^T \mathbf{w}) \mathbf{v}_i.$$

具体算法如下.

算法 2.14 (重正交化的 Lanczos 方法) 给定矩阵 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 满足 $\mathbf{A}^T = \mathbf{A}$, 向量 $\mathbf{v}_1 \in \mathbb{R}^n$ 满足 $\|\mathbf{v}_1\|_2 = 1$ 以及正整数 k , 本算法计算一个长度为 k 的 Lanczos 分解: $\mathbf{A}\mathbf{V}_k = \mathbf{V}_k \mathbf{T}_k + \beta_k \mathbf{v}_{k+1} \mathbf{e}_k^T$.

$\beta_0 = 0; \mathbf{v}_0 = \mathbf{0};$

for $j = 1 : k$

$\mathbf{w} = \mathbf{A}\mathbf{v}_j; \alpha_j = \mathbf{v}_j^T \mathbf{w};$

$\mathbf{w} = \mathbf{w} - \alpha_j \mathbf{v}_j - \beta_{j-1} \mathbf{v}_{j-1};$

$\mathbf{w} = \mathbf{w} - \sum_{i=1}^{j-1} (\mathbf{v}_i^T \mathbf{w}) \mathbf{v}_i;$

$\beta_j = \|\mathbf{w}\|_2;$

if $\beta_j = 0$

stop

else

$\mathbf{v}_{j+1} = \mathbf{w} / \beta_j;$

end

end

当然, 在具体计算时, 也可以用修正的 Gram-Schmidt 正交化方法来实现上式的计算. 有时为了使算法所产生的 \mathbf{v}_j 之间有更好的正交性, 也可以连续进行两次重正交化, 即所谓的完全重正交化.

下面再给出一个具体的例子, 观察一下重正交化的 Lanczos 过程的效果.

例 2.7 考虑矩阵

$$\mathbf{A} = \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & 2 & 2 & \cdots & 2 \\ 1 & 2 & 3 & \cdots & 3 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 2 & 3 & \cdots & n \end{bmatrix}. \quad (2.37)$$

随机地选取一个初始向量 \mathbf{v}_1 , 并且取 $n = 1000, k = 50$. 若用算法 2.13, 则计算得到的矩阵 \mathbf{V}_{50} 满足

$$\|\mathbf{V}_{50}^T \mathbf{V}_{50} - \mathbf{I}_{50}\|_F = 9.5938.$$

而若用重正交化的 Lanczos 算法, 则计算得到的矩阵 V_{50} 满足

$$\|V_{50}^T V_{50} - I_{50}\|_F = 2.0073 \times 10^{-13}.$$

由此观察到, 重正交化的 Lanczos 过程的效果是明显的.

2.5 投影方法

大多数已有的求解线性方程组迭代方法都按某种方式使用一种投影过程. 投影过程表达了从一个子空间推断出线性方程组近似解的一种规范方式. 本节在一般的框架下描述这种规范方式, 并给出一些理论. 进一步, 对一维情形的投影方法给出了较为详细的描述, 这给第 4 章中更为复杂的子空间投影方法提供了一个预览.

2.5.1 投影算子及其性质

投影算子 P 是从 \mathbb{C}^n 到其自身的一个幂等线性映射, 即满足

$$P^2 = P.$$

由此定义可得下述一些简单的性质. 首先, 如果 P 是一个投影, 则 $I - P$ 也是投影, 且下述关系成立,

$$\mathcal{N}(P) = \mathcal{R}(I - P).$$

另外, 两个子空间 $\mathcal{N}(P)$ 和 $\mathcal{R}(P)$ 的交只有元素 $\mathbf{0}$, 即 $\mathcal{N}(P) \cap \mathcal{R}(P) = \{\mathbf{0}\}$. 事实上, 如果向量 x 属于 $\mathcal{R}(P)$, 则由幂等性有 $Px = x$. 如果它也属于 $\mathcal{N}(P)$, 则 $Px = \mathbf{0}$, 因此 $x = Px = \mathbf{0}$. 此外, \mathbb{C}^n 的每个元素可被写成 $x = Px + (I - P)x$, 因此, 空间 \mathbb{C}^n 可被分解成直和

$$\mathbb{C}^n = \mathcal{N}(P) \oplus \mathcal{R}(P).$$

反之, 每个形成 \mathbb{C}^n 的直和的子空间对 \mathcal{K} 和 \mathcal{S} 唯一地定义一个投影算子, 使得 $\mathcal{R}(P) = \mathcal{K}$ 和 $\mathcal{N}(P) = \mathcal{S}$. 相应的投影算子 P 映 \mathbb{C}^n 的一个元素 x 为分量 x_1 , 其中 x_1 是相应直和唯一分解 $x = x_1 + x_2$ 中的 \mathcal{K} 分量.

事实上, 这种对应是唯一的, 即任意投影算子 P 可由给定的两个子空间完全确定: P 的值域 \mathcal{K} 和它的核 \mathcal{S} (亦即 $I - P$ 的值域). 若对任意的 x , 向量 Px 满足

$$Px \in \mathcal{K}, \quad x - Px \in \mathcal{S},$$

则称线性映射 P 沿 (或平行于) 子空间 \mathcal{S} 将 x 投影到 \mathcal{K} 上. 如果 P 的秩为 m , 则 $I - P$ 的值域是 $n - m$ 维的. 因此, 通过其 m 维的正交补 $\mathcal{L} = \mathcal{S}^\perp$ 来定义 \mathcal{S} 是自然的. 对任意的 x , 上述定义 $u = Px$ 的条件变为

$$u \in \mathcal{K}, \tag{2.38}$$

$$x - u \perp \mathcal{L}. \tag{2.39}$$

这定义了一个正交于子空间 \mathcal{L} 的到 \mathcal{K} 上的投影算子 P . 式 (2.38) 建立了 m 个自由度, 式 (2.39) 给出了从这些自由度来定义 Px 的 m 个约束.

现在的问题是: 任给两个都是 m 维的子空间 \mathcal{K} 和 \mathcal{L} , 是否总能通过式 (2.38) 和式 (2.39) 定义一个正交于 \mathcal{L} 的到 \mathcal{K} 上的投影算子 P ? 下面的引理回答这个问题.

引理 2.1 给定两个 m 维子空间 \mathcal{K} 和 \mathcal{L} , 下面两个条件是等价的:

- (1) \mathcal{K} 中没有非零向量正交于 \mathcal{L} .
- (2) 对任意的 $x \in \mathbb{C}^n$, 存在唯一向量 u 满足式 (2.38) 和 (2.39).

引理 2.1 说明给定两个满足条件 $\mathcal{K} \cap \mathcal{L}^\perp = \{0\}$ 的子空间 \mathcal{K} 和 \mathcal{L} , 存在一个正交于 \mathcal{L} 的到 \mathcal{K} 上的投影算子 P , 它用式 (2.38) 和 (2.39) 定义任意向量 x 的投影向量 u . 此投影算子满足

$$\mathcal{R}(P) = \mathcal{K}, \quad \mathcal{N}(P) = \mathcal{L}^\perp.$$

特别地, 条件 $Px = 0$ 转化为 $x \in \mathcal{N}(P)$, 而这意味着 $x \in \mathcal{L}^\perp$. 反之亦真. 这可得下面的性质:

$$Px = 0 \iff x \in \mathcal{L}^\perp.$$

为了得到一般投影算子的矩阵表示, 需要两组基: 子空间 $\mathcal{K} = \mathcal{R}(P)$ 的一组基 $V = [v_1, \dots, v_m]$ 和子空间 \mathcal{L} 的一组基 $W = [w_1, \dots, w_m]$. 当

$$(v_i, w_j) = \delta_{ij} = \begin{cases} 1, & \text{若 } i = j \\ 0, & \text{若 } i \neq j \end{cases},$$

时, 称这两组基是双正交的, 写成矩阵形式为 $W^H V = I$. 由于 Px 属于 \mathcal{K} , 故可表示为 $Px = Vy$. 约束 $x - Px \perp \mathcal{L}$ 等价于条件

$$(x - Vy, w_i) = 0, \quad i = 1, \dots, m.$$

写成矩阵形式, 即

$$W^H(x - Vy) = 0. \quad (2.40)$$

如果两组基是双正交的, 则得 $y = W^H x$. 因此, 此时 $Px = VW^H x$, 这导致 P 的矩阵表示

$$P = VW^H. \quad (2.41)$$

对两组基 V 和 W 不是双正交的情形, 由式 (2.40) 易见

$$P = V(W^H V)^{-1} W^H.$$

如果假设 \mathcal{K} 中没有向量正交于 \mathcal{L} , 则可以证明 $m \times m$ 阶矩阵 $W^H V$ 是非奇异的.

当子空间 $\mathcal{L} = \mathcal{K}$ 时, 称投影算子 P 为到 \mathcal{K} 上的正交投影算子. 非正交投影算子称为斜投影算子. 对于正交投影算子, 有

$$\mathcal{N}(P) = \mathcal{R}(P)^\perp.$$

因此, 可由下面对任意向量 x 都要满足的要求来定义一个正交投影算子:

$$Px \in \mathcal{K} \text{ 且 } (I - P)x \perp \mathcal{K},$$

或等价地, 有

$$Px \in \mathcal{K} \text{ 且 } ((I-P)x, y) = 0, \quad \forall y \in \mathcal{K}.$$

考虑映射 P 的伴随映射 P^H , 它定义为

$$(P^H x, y) = (x, Py), \quad \forall x, y. \quad (2.42)$$

首先注意到 P^H 也是一个投影算子, 因为对所有的 x 和 y ,

$$((P^H)^2 x, y) = (P^H x, Py) = (x, P^2 y) = (x, Py) = (P^H x, y).$$

由式 (2.42), 得

$$\mathcal{N}(P^H) = \mathcal{R}(P)^\perp, \quad (2.43)$$

$$\mathcal{N}(P) = \mathcal{R}(P^H)^\perp. \quad (2.44)$$

上述关系导致下面的性质.

性质 2.1 投影算子是正交的当且仅当它是 Hermite 的.

任给一个矩阵 $V \in \mathbb{C}^{n \times m}$, 其列满足 $v_i^H v_i = 1 (i = 1, 2, \dots, m)$. 则其列形成 $\mathcal{K} = \mathcal{R}(P)$ 的一组标准正交基, 可用矩阵 $P = VV^H$ 表示 P . 这是投影算子矩阵表示 (2.41) 的一种特殊情况. 除了是幂等的, 与此矩阵相应的线性映射满足上述给出的特征, 即

$$VV^H x \in \mathcal{K} \text{ 且 } (I - VV^H)x \in \mathcal{K}^\perp.$$

重要的是注意到正交投影算子 P 的这种表示不是唯一的. 事实上, 任意标准正交基 V 将给出 P 的一种不同的形式的表示. 因此, 对 \mathcal{K} 的任意两组正交基 V_1 和 V_2 , 必须有 $V_1 V_1^H = V_2 V_2^H$, 这一等式可独立地证明.

下面来看看正交投影的性质. 当 P 是一个正交投影算子时, 分解 $x = Px + (I-P)x$ 中的两个向量是正交的, 得到下述关系:

$$\|x\|_2^2 = \|Px\|_2^2 + \|(I-P)x\|_2^2.$$

因此, 对任意 x , 有

$$\|Px\|_2 \leq \|x\|_2.$$

故对所有 $x \in \mathbb{C}^n$, $\|Px\|_2 / \|x\|_2$ 的最大值不超过 1, 另外, 对 $\mathcal{R}(P)$ 中任意元素达到 1. 所以对任意正交投影算子 P , 有

$$\|P\|_2 = 1.$$

正交投影算子只有两个特征值: 0 和 1. P 的值域中的任意向量都是相应于特征值 1 的特征向量, 核中的任意向量都是相应于特征值 0 的特征向量.

下面建立关于正交投影算子的一个重要的最优性性质.

定理 2.13 设 P 是到子空间 \mathcal{K} 上的正交投影算子, 则对任意给定的向量 $x \in \mathbb{C}^n$, 下式成立:

$$\min_{y \in \mathcal{K}} \|x - y\|_2 = \|x - Px\|_2. \quad (2.45)$$

证明 设 y 为 \mathcal{K} 中任意向量, 考虑它到 x 的距离的平方. 由于 $x - Px$ 正交于 \mathcal{K} , 而 $Px - y$ 属于 \mathcal{K} , 则

$$\|x - y\|_2^2 = \|x - Px + (Px - y)\|_2^2 = \|x - Px\|_2^2 + \|Px - y\|_2^2.$$

因此, 对任意 $y \in \mathcal{K}$, $\|x - y\|_2^2 \geq \|x - Px\|_2^2$. 注意到对 $y = Px$ 达到极小, 这就证明了式 (2.45). 证毕. \square

通过对子空间 \mathcal{K} 上的正交投影算子 P 定义 $y^* \equiv Px$, 可将上述结果重新表示成一个充分必要条件的形式, 这个条件可以确定一个向量 x 在最小二乘意义下的最佳逼近.

推论 2.2 设给定一个子空间 \mathcal{K} 和一个向量 $x \in \mathbb{C}^n$, 则

$$\min_{y \in \mathcal{K}} \|x - y\|_2 = \|x - y^*\|_2,$$

当且仅当满足下述条件:

$$\begin{cases} y^* \in \mathcal{K}, \\ x - y^* \perp \mathcal{K}. \end{cases}$$

2.5.2 投影方法的基本框架

考虑线性方程组 $Ax = b$, 其中 $A \in \mathbb{R}^{n \times n}$. 这里, 符号 A 既用来表示矩阵, 也用来表示它所代表的 \mathbb{R}^n 中的线性映射. 投影方法的思想是从 \mathbb{R}^n 的一个子空间推断出问题的一个近似解. 如果 \mathcal{K} 是这种候选近似子空间, 或称“搜索子空间”, 并且维数为 m , 则一般地为推断出这种近似必须加 m 个约束. 描述这些约束的一个典型的方式是强加 m 个正交性条件. 特别地, 残差向量 $b - Ax$ 限制为与 m 个线性无关的向量正交. 这就定义了另一个 m 维子空间 \mathcal{L} , 称其为“约束子空间”. 这种简单的框架称为 Petrov-Galerkin 条件, 它在许多不同的数学方法中是很普遍的.

投影方法有两大类: 正交投影方法和斜投影方法. 在正交投影技术中, 子空间 \mathcal{L} 与 \mathcal{K} 是相同的. 在斜投影方法中 \mathcal{L} 与 \mathcal{K} 不相同, 甚至彼此完全无关. 这种差别非常重要并由此给出不同类型的算法.

设 $A \in \mathbb{R}^{n \times n}$ 且 \mathcal{K} 和 \mathcal{L} 为 \mathbb{R}^n 的两个 m 维子空间. 正交于 \mathcal{L} 的到 \mathcal{K} 上的投影技术是一个过程: 通过条件 $\tilde{x} \in \mathcal{K}$ 并且新的残差向量正交于 \mathcal{L} 来求 $Ax = b$ 的一个近似解 \tilde{x} , 即

$$\text{求 } \tilde{x} \in \mathcal{K}, \text{ 使得 } b - A\tilde{x} \perp \mathcal{L}.$$

如果希望利用解的初始值 $x^{(0)}$ 的信息, 则必须在仿射空间 $x^{(0)} + \mathcal{K}$ 中而不是在齐次的向量空间 \mathcal{K} 中寻求近似解. 这需要对上述公式进行稍微的修改. 近似解问题需要重新定义为

$$\tilde{x} \in x^{(0)} + \mathcal{K}, \text{ 使得 } b - A\tilde{x} \perp \mathcal{L}.$$

注意到, 如果将 \tilde{x} 写成形式 $\tilde{x} = x^{(0)} + z$, 且初始残差向量 $r^{(0)}$ 定义为 $r^{(0)} = b - Ax^{(0)}$, 则上述方程变为

$$b - A(x^{(0)} + z) \perp \mathcal{L} \text{ 或 } r^{(0)} - Az \perp \mathcal{L}.$$

换言之, 可定义近似解为

$$\tilde{x} = x^{(0)} + z, \quad z \in \mathcal{K}, \quad (2.46)$$

$$(r^{(0)} - Az, w) = 0, \quad \forall w \in \mathcal{L}. \quad (2.47)$$

这是最一般形式的一个基本投影步. 大多数标准技术都连续地使用这种投影. 投影方法为科学计算中许多著名的方法提供了统一的框架. 事实上, 几乎所有的基本迭代法都可视为投影技术. 只要通过 m 个自由度 (子空间 \mathcal{K}) 和 m 个约束 (子空间 \mathcal{L}) 来定义一个近似, 就可得到一个投影过程.

正交投影方法相应于两个子空间 \mathcal{K} 和 \mathcal{L} 相等的特殊情形. 这种差别在 Heimate 的情形下特别重要, 因为在此情形下, 需要保证被投影的问题是 Heimate 的. 当 $\mathcal{K} = \mathcal{L}$ 时, Petrov-Galerkin 条件也简称为 Galerkin 条件.

设 $V = [v_1, \dots, v_m]$ 是一个 $n \times m$ 阶矩阵, 其列向量形成 \mathcal{K} 的一组基. 类似地, $W = [w_1, \dots, w_m]$ 也是一个 $n \times m$ 阶矩阵, 其列向量形成 \mathcal{L} 的一组基. 如果将近似解写成

$$x = x^{(0)} + Vy,$$

则正交性条件立即导致下面关于向量 y 的方程组:

$$W^T AVy = W^T r^{(0)}.$$

如果假设 $m \times m$ 阶矩阵 $W^T AV$ 是非奇异的, 则有下列关于近似解 x 的表达式:

$$x = x^{(0)} + V(W^T AV)^{-1} W^T r^{(0)}. \quad (2.48)$$

在许多算法中, 没有必要形成矩阵 $W^T AV$, 因此它可作为对算法的一个“副产品”供使用. 典型投影方法由下述算法表示.

算法 2.15 (投影方法的基本框架)

- 步 1, 选取一对子空间 \mathcal{K} 和 \mathcal{L} .
- 步 2, 选取 \mathcal{K} 和 \mathcal{L} 的基 $V = [v_1, \dots, v_m]$ 和 $W = [w_1, \dots, w_m]$.
- 步 3, 计算残差 $r = b - Ax$.
- 步 4, 解方程组 $W^T AVy = W^T r$ 得到 y .
- 步 5, 计算近似解 $x := x + Vy$.

只有当矩阵 $W^T AV$ 是非奇异的, 近似解才是有意义的. 即使当 A 是非奇异的, 也不能保证 $W^T AV$ 是非奇异的. 看下面的例子.

例 2.8 考虑矩阵

$$A = \begin{bmatrix} O & I \\ I & I \end{bmatrix},$$

式中: I 为 m 阶单位矩阵且 O 为 m 阶零矩阵. 设 $V = W = [e_1, e_2, \dots, e_m]$, 其中 e_i 是第 i 个分量为 1 而其余分量均为 0 的 $2m$ 维列向量. 尽管 A 是非奇异的, 而矩阵 $W^T AV = O$, 因此是奇异的.

下面的命题表明, 有两种重要的特殊情形可以保证 $W^T A V$ 的非奇异性.

命题 2.1 设 A , K 和 L 满足下面两个条件之一:

(1) A 是正定的且 $L = K$.

(2) A 是非奇异的且 $L = AK$.

则矩阵 $B = W^T A V$ 对 K 和 L 的任意基 V 和 W 均是非奇异的.

证明 (1) 设 V 和 W 分别为 K 和 L 的任意一组基. 事实上, 由于 $L = K$, W 总可以表示为 $W = VG$, 其中 G 是一个非奇异 m 阶矩阵, 则

$$B = W^T A V = G^T V^T A V.$$

由于 A 是正定的, 所以 $V^T A V$ 正定, 从而证明了 B 是非奇异的.

(2) 设 V 和 W 分别为 K 和 L 的任意一组基. 由于 $L = AK$, 此时 W 总可表示为 $W = AVG$, 其中 G 是一个非奇异 m 阶矩阵, 则

$$B = W^T A V = G^T (AV)^T A V.$$

由于 A 是非奇异的, $n \times m$ 阶矩阵 AV 是满秩的, 因此 $(AV)^T A V$ 是非奇异的. 由此及上式证明了 B 是非奇异的. 证毕. \square

注 2.14 考虑特殊情形: A 是实对称的且使用正交投影, 此时, 对 L 和 K 使用相同的基, 因为它们有相同的子空间, 则投影矩阵 $B = V^T A V$ 是对称的. 另外, 如果矩阵 A 是对称正定的, 则 B 也是对称正定的.

进一步, 分别考察命题 2.1 中的两个条件所对应的最优性结果.

命题 2.2 假设 A 是对称正定的且 $L = K$, 初始向量为 $x^{(0)}$. 则向量 \tilde{x} 为到 K 上的一个正交投影方法的结果, 当且仅当 \tilde{x} 在 $x^{(0)} + K$ 上极小化误差的 A 范数 (能量范数), 即当且仅当

$$\mathcal{E}(\tilde{x}) = \min_{x \in x^{(0)} + K} \mathcal{E}(x),$$

式中: $\mathcal{E}(x) \equiv \|x^* - x\|_A = \sqrt{(A(x^* - x), x^* - x)}$.

证明 由推论 2.2, 为使 \tilde{x} 极小化 $\mathcal{E}(x)$, 当且仅当 $x^* - \tilde{x}$ 与 K 的所有向量均 A 正交. 由此得

$$(A(x^* - \tilde{x}), v) = 0, \quad \forall v \in K,$$

或等价地, 有

$$(b - A\tilde{x}, v) = 0, \quad \forall v \in K,$$

这正是对近似解 \tilde{x} 定义一个正交投影的 Galerkin 条件. 证毕. \square

命题 2.3 设 A 是任意方阵且 $L = AK$, 初始向量为 $x^{(0)}$. 则向量 \tilde{x} 为正交于 L 的到 K 上一个斜投影方法的结果, 当且仅当 \tilde{x} 在 $x \in x^{(0)} + K$ 上极小化残差向量 $r(x) = b - Ax$ 的 2-范数, 即当且仅当

$$\|r(\tilde{x})\|_2 = \min_{x \in x^{(0)} + K} \|r(x)\|_2.$$

证明 如前所见, 为使 \tilde{x} 极小化 $\|r(x)\|_2$, 其充分必要条件是 $b - A\tilde{x}$ 与所有形如 $v = Ay$ 的向量均正交, 其中 $y \in K$. 亦即

$$(b - A\tilde{x}, v) = 0, \quad \forall v \in AK,$$

这正是定义近似解 \tilde{x} 的 Petrov-Galerkin 条件. 证毕. \square

值得指出的是, 上述命题中 A 不需要是非奇异的. 当 A 为奇异时, 可能有无穷多个向量 \tilde{x} 满足最优性条件.

下面给出用投影算子术语的解释. 对 $\mathcal{L} = K$ 和 $\mathcal{L} = AK$ 的情形, 可容易地用在初始残差或初始误差上正交投影算子的作用来解释投影方法的结果. 先考虑 $\mathcal{L} = AK$ 的情形, 设初始残差为 $r^{(0)} = b - Ax^{(0)}$, 且 $\tilde{r} = b - A\tilde{x}$ 为投影过程之后所得的残差, 则

$$\tilde{r} = b - A(x^{(0)} + z) = r^{(0)} - Az.$$

其中 z 由强加条件 $r^{(0)} - Az \perp AK$ 而得到. 因此, 向量 Az 是向量 $r^{(0)}$ 到子空间 AK 上的正交投影. 从而有下面的命题.

命题 2.4 设 \tilde{x} 为由正交于 $\mathcal{L} = AK$ 的到 K 上的一个投影过程所得的近似解, 且设 $\tilde{r} = b - A\tilde{x}$ 为相应的残差, 则

$$\tilde{r} = (I - P)r^{(0)}, \quad (2.49)$$

式中: P 为子空间 AK 上的正交投影算子.

命题 2.4 说明一个投影步后所得残差的 2-范数将不超过初始残差的 2-范数, 即

$$\|\tilde{r}\|_2 \leq \|r^{(0)}\|_2,$$

这是一个已经建立了的结果. 这类方法称为残差投影方法.

现在考虑 $\mathcal{L} = K$ 且 A 是对称正定的情形. 设 $e^{(0)} = x^* - x^{(0)}$ 为初始误差, 类似地, 设 $\tilde{e} = x^* - \tilde{x}$, 其中 $\tilde{x} = x^{(0)} + z$ 为由投影步所得的近似解. 由式 (2.49), 得

$$A\tilde{e} = A(x^* - \tilde{x}) = \tilde{r} = A(e^{(0)} - z) = r^{(0)} - Az,$$

式中: z 是由限制残差向量 $r^{(0)} - Az \perp K$ 而得到, 即

$$(r^{(0)} - Az, w) = 0, \quad \forall w \in K \iff (A(e^{(0)} - z), w) = 0, \quad \forall w \in K.$$

由于 A 是对称正定的, 它定义了一种通常表示为 $(\cdot, \cdot)_A$ 的内积, 则条件变为

$$(e^{(0)} - z, w)_A = 0, \quad \forall w \in K.$$

此条件可解释为: 向量 z 是初始误差向量 $e^{(0)}$ 到子空间 K 上的 A 正交投影.

命题 2.5 设 \tilde{x} 为到 K 上的一个正交投影过程所得的近似解, 且设 $\tilde{e} = x^* - \tilde{x}$ 为相应的误差, 则

$$\tilde{e} = (I - P_A)e^{(0)},$$

式中: P_A 为子空间 K 上的投影算子, 它关于 A 内积是正交的.

命题 2.5 表明一个投影步后所得的误差向量的 A 范数不超过初始误差的 A 范数, 即

$$\|\tilde{e}\|_A \leq \|e^{(0)}\|_A,$$

这正是所期望的, 因为已知误差的 A 范数在 $x^{(0)} + \mathcal{K}$ 中被极小化. 这类方法称为误差投影方法.

2.5.3 一维投影方法

这里给出几种基于一维投影过程的投影方法. 以下用 r 表示当前近似解 x 的残差向量 $r = b - Ax$. 为避免下标, 用 “:=” 表示向量校正, 即 $x := x + \alpha r$ 意味着 “计算 $x + \alpha r$ 并用结果覆盖当前的 x ”.

一维投影过程可由下面定义, 即

$$\mathcal{K} = \text{span}\{v\} \text{ 且 } \mathcal{L} = \text{span}\{w\},$$

式中: v 和 w 为两个向量.

此时, 新的近似取为 $x := x + \alpha v$ 且 Petrov-Galerkin 条件 $r - Az \perp w$ 导致

$$\alpha = \frac{(r, w)}{(Av, w)}.$$

下面考虑三种流行的选择.

1) 最速下降法

最速下降法对矩阵 A 为对称正定的情形来定义. 它由每一步取 $v = r$ 和 $w = r$ 组成 (即一维正交投影). 这导致下面描述的一个迭代.

算法 2.16 (最速下降法)

步 1, $r := b - Ax$.

步 2, $\alpha \leftarrow (r, r)/(Ar, r)$.

步 3, $x := x + \alpha r$.

此迭代的每一步在所有形式为 $x + \alpha r$ 的向量上极小化

$$f(x) = \|x - x^*\|_A^2 = (A(x - x^*), x - x^*),$$

式中: r 为负梯度方向 $-\nabla f(x)$.

负梯度方向是使 $f(x)$ 局部最快速下降的方向. 其次, 证明当 A 是对称正定的时, 收敛性是有保证的. 它是下面著名的 Kantorovich (康托洛维奇) 不等式的一个结论.

引理 2.2 (Kantorovich 不等式) 设 B 为任意对称正定实矩阵且 $\lambda_{\max}, \lambda_{\min}$ 为其最大、最小特征值, 则

$$\frac{(Bx, x)(B^{-1}x, x)}{(x, x)^2} \leq \frac{(\lambda_{\max} + \lambda_{\min})^2}{4\lambda_{\max}\lambda_{\min}}, \quad \forall x \neq 0. \quad (2.50)$$

定理 2.14 设 A 是对称正定的. 则由算法 2.16 生成的误差向量 $e^{(k)} = x^* - x^{(k)}$ 的 A 范数满足关系

$$\|e^{(k+1)}\|_A \leq \frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\max} + \lambda_{\min}} \|e^{(k)}\|_A, \quad (2.51)$$

且算法 2.16 对任意初始向量 $x^{(0)}$ 都是收敛的.

证明 首先注意到 $\|e^{(k+1)}\|_A^2 = (Ae^{(k+1)}, e^{(k+1)}) = (r^{(k+1)}, e^{(k+1)})$, 则由简单的代入有

$$\|e^{(k+1)}\|_A^2 = (r^{(k+1)}, e^{(k)} - \alpha_k r^{(k)}).$$

由于新残差向量 $r^{(k+1)}$ 必须正交于搜索方向 $r^{(k)}$, 上式右端项中第 2 项为零. 因此

$$\begin{aligned} \|e^{(k+1)}\|_A^2 &= (r^{(k)} - \alpha_k Ar^{(k)}, e^{(k)}) \\ &= (r^{(k)}, A^{-1}r^{(k)}) - \alpha_k (r^{(k)}, r^{(k)}) \\ &= \|e^{(k)}\|_A^2 \left(1 - \frac{(r^{(k)}, r^{(k)})}{(r^{(k)}, Ar^{(k)})} \cdot \frac{(r^{(k)}, r^{(k)})}{(r^{(k)}, A^{-1}r^{(k)})} \right). \end{aligned}$$

应用 Kantorovich 不等式 (2.50) 即得定理的结论. 证毕. \square

2) 极小残差 (MR) 法

现在假设 A 不必为对称的而只是正定的, 即其对称部分 $A + A^T$ 是对称正定的. 在每一步取 $v = r$ 和 $w = Ar$, 可得下述迭代过程.

算法 2.17 (极小残差法)

- 步 1, $r := b - Ax$.
- 步 2, $\alpha := (Ar, r)/(Ar, Ar)$.
- 步 3, $x := x + \alpha r$.

这里, 每一步在方向 r 上极小化 $\|b - Ax\|_2^2$. 如下定理所述, 迭代在 A 为正定时收敛.

定理 2.15 设 A 为实正定矩阵, 且设

$$\mu = \lambda_{\min}(A + A^T)/2, \quad \sigma = \|A\|_2.$$

则由算法 2.17 产生的残差向量满足关系

$$\|r^{(k+1)}\|_2 \leq \left(1 - \frac{\mu^2}{\sigma^2} \right)^{1/2} \|r^{(k)}\|_2, \quad (2.52)$$

且算法 2.17 对任意初始向量 $x^{(0)}$ 均收敛.

证明 类似于最速下降法来处理, 有关系式

$$\|r^{(k+1)}\|_2^2 = (r^{(k)} - \alpha_k Ar^{(k)}, r^{(k)} - \alpha_k Ar^{(k)})$$

$$= (\mathbf{r}^{(k)} - \alpha_k \mathbf{A} \mathbf{r}^{(k)}, \mathbf{r}^{(k)}) - \alpha_k (\mathbf{r}^{(k)} - \alpha_k \mathbf{A} \mathbf{r}^{(k)}, \mathbf{A} \mathbf{r}^{(k)}).$$

由构造, 新的残差向量 $\mathbf{r}^{(k)} - \alpha_k \mathbf{A} \mathbf{r}^{(k)}$ 必须正交于搜索方向 $\mathbf{A} \mathbf{r}^{(k)}$, 因此, 上式右端项中第 2 项为零, 且得

$$\begin{aligned} \|\mathbf{r}^{(k+1)}\|_2^2 &= (\mathbf{r}^{(k)} - \alpha_k \mathbf{A} \mathbf{r}^{(k)}, \mathbf{r}^{(k)}) = (\mathbf{r}^{(k)}, \mathbf{r}^{(k)}) - \alpha_k (\mathbf{A} \mathbf{r}^{(k)}, \mathbf{r}^{(k)}) \\ &= \|\mathbf{r}^{(k)}\|_2^2 \left(1 - \frac{(\mathbf{A} \mathbf{r}^{(k)}, \mathbf{r}^{(k)})}{(\mathbf{r}^{(k)}, \mathbf{r}^{(k)})} \cdot \frac{(\mathbf{A} \mathbf{r}^{(k)}, \mathbf{r}^{(k)})}{(\mathbf{A} \mathbf{r}^{(k)}, \mathbf{A} \mathbf{r}^{(k)})} \right) \\ &= \|\mathbf{r}^{(k)}\|_2^2 \left(1 - \frac{(\mathbf{A} \mathbf{r}^{(k)}, \mathbf{r}^{(k)})^2}{(\mathbf{r}^{(k)}, \mathbf{r}^{(k)})^2} \cdot \frac{\|\mathbf{r}^{(k)}\|_2^2}{\|\mathbf{A} \mathbf{r}^{(k)}\|_2^2} \right). \end{aligned} \quad (2.53)$$

因为 \mathbf{A} 为实正定矩阵, 从而有

$$\frac{(\mathbf{A} \mathbf{x}, \mathbf{x})}{(\mathbf{x}, \mathbf{x})} = \frac{1}{2} \frac{((\mathbf{A} + \mathbf{A}^T) \mathbf{x}, \mathbf{x})}{(\mathbf{x}, \mathbf{x})} \geq \mu > 0, \quad (2.54)$$

式中: $\mu = \lambda_{\min}(\mathbf{A} + \mathbf{A}^T)/2$.

由不等式 $\|\mathbf{A} \mathbf{r}^{(k)}\|_2 \leq \|\mathbf{A}\|_2 \|\mathbf{r}^{(k)}\|_2$ 立即可得定理的结论. 证毕. \square

3) 残差范数最速下降法

在残差范数最速下降法中, 只需 \mathbf{A} 为非奇异矩阵. 在每一步, 算法使用 $\mathbf{v} = \mathbf{A}^T \mathbf{r}$ 和 $\mathbf{w} = \mathbf{A} \mathbf{v}$, 给出下面运算序列:

$$\mathbf{r} := \mathbf{b} - \mathbf{A} \mathbf{x}, \quad \mathbf{v} = \mathbf{A}^T \mathbf{r},$$

$$\alpha := \|\mathbf{v}\|_2^2 / \|\mathbf{A} \mathbf{v}\|_2^2,$$

$$\mathbf{x} := \mathbf{x} + \alpha \mathbf{v}.$$

注意到基于上述运算序列的算法需要三个矩阵向量乘积, 这是本节其他算法的 3 倍. 通过不同的方式计算残差, 每步中矩阵向量乘积的个数可降为两个. 这种变形如下:

算法 2.18 (残差范数最速下降)

步 1, 计算 $\mathbf{r} := \mathbf{b} - \mathbf{A} \mathbf{x}$.

步 2, $\mathbf{v} := \mathbf{A}^T \mathbf{r}$.

步 3, 计算 $\mathbf{A} \mathbf{v}$ 和 $\alpha := \|\mathbf{v}\|_2^2 / \|\mathbf{A} \mathbf{v}\|_2^2$.

步 4, $\mathbf{x} := \mathbf{x} + \alpha \mathbf{v}$.

步 5, $\mathbf{r} := \mathbf{r} - \alpha \mathbf{A} \mathbf{v}$.

这里, 每一步在方向 $-\nabla f(\mathbf{x})$ 上极小化 $\|\mathbf{b} - \mathbf{A} \mathbf{x}\|_2^2$. 由此得出, 这等价于最速下降算法应用于法方程. 因为当 \mathbf{A} 是非奇异时, $\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}$ 是正定的, 则根据定理 2.15, 只要 \mathbf{A} 是非奇异的, 方法将收敛.

习题 2

2.1 设向量 $\mathbf{x} = (2, 4, 0, 5, 1, 3)^T$. 求一个 Householder 变换 \mathbf{H} 和一个正常数 α , 使得 $\mathbf{H}\mathbf{x} = (2, 4, \alpha, 5, 0, 0)^T$.

2.2 假定 \mathbf{x} 和 \mathbf{y} 是 \mathbb{R}^n 中满足 $\|\mathbf{x}\|_2 = \|\mathbf{y}\|_2 = 1$ 的两个向量. 给出一种使用 Givens 变换的算法, 构造一个正交矩阵 \mathbf{Q} , 使得 $\mathbf{Q}\mathbf{x} = \mathbf{y}$.

2.3 假定 \mathbf{x} 和 \mathbf{y} 是 \mathbb{R}^n 中的两个非零向量. 给出一种算法来构造一个 Householder 矩阵 \mathbf{H} , 使得 $\mathbf{H}\mathbf{x} = \alpha\mathbf{y}$, 其中 $\alpha \in \mathbb{R}$.

2.4 设矩阵

$$\mathbf{A} = \begin{bmatrix} 0 & 4 & 1 \\ 1 & 1 & 1 \\ 0 & 3 & 2 \end{bmatrix},$$

利用 Householder 变换求 \mathbf{A} 的 QR 分解.

2.5 设 $\mathbf{H} \in \mathbb{R}^{n \times n}$ 是上 Hessenberg 矩阵, $\mathbf{R} \in \mathbb{R}^{n \times n}$ 是非奇异的上三角矩阵. 证明: $\mathbf{R}\mathbf{H}\mathbf{R}^{-1}$ 仍是上 Hessenberg 矩阵.

2.6 设

$$\mathbf{A} = \begin{bmatrix} 2 & 2 & 1 \\ 0 & 2 & 2 \\ 2 & 1 & 2 \end{bmatrix},$$

利用 Givens 变换求矩阵 \mathbf{A} 的 QR 分解.

2.7 设矩阵 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 可以分解为 $\mathbf{A} = \mathbf{Q}\mathbf{R}$, 其中 \mathbf{Q} 为正交矩阵, \mathbf{R} 是对角元全为正的上三角矩阵. 记 \mathbf{A} 的第 i 列为 \mathbf{a}_i . 证明:

$$|\det(\mathbf{A})| \leq \prod_{i=1}^n \|\mathbf{a}_i\|_2.$$

2.8 设 \mathbf{A} 为 $n \times n$ 非奇异实矩阵, 其 QR 分解为 $\mathbf{A} = \mathbf{Q}\mathbf{R}$. 记 $\mathbf{B} = \mathbf{R}\mathbf{Q}$, 证明:

(1) 若 \mathbf{A} 是对称的, 则 \mathbf{B} 也是对称的;

(2) 若 \mathbf{A} 是上 Hessenberg 矩阵, 则 \mathbf{B} 也是上 Hessenberg 矩阵.

2.9 假设 $\mathbf{v} = p(\mathbf{A})\mathbf{u}$, 其中 $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$, p 是多项式. 证明: 对每个 $k = 1, 2, \dots, n$, 都有 $p(\mathbf{A})\mathcal{K}_k(\mathbf{A}, \mathbf{u}) = \mathcal{K}_k(\mathbf{A}, \mathbf{v})$.

2.10 设 $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$ 是矩阵 \mathbf{A} 的 k 个线性无关的特征向量, 令

$$\mathbf{v} = \alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2 + \dots + \alpha_k \mathbf{x}_k,$$

其中系数 $\alpha_1, \alpha_2, \dots, \alpha_k$ 都是非零的. 证明:

$$\mathcal{K}_k(\mathbf{A}, \mathbf{v}) = \text{span}\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\}.$$

2.11 证明: 任一正交投影算子 \mathbf{P} 均可表示为 $\mathbf{P} = \mathbf{V}\mathbf{V}^H$, 其中 \mathbf{V} 满足 $\mathbf{V}^H\mathbf{V} = \mathbf{I}$.

2.12 设 A 和 B 是从 \mathbb{C}^n 到其自身的两个投影算子. 证明:

(1) $A + B$ 是投影算子当且仅当 $BA = AB = O$;

(2) 若 $AB = BA$, 则 $A + B - AB$ 也是投影算子.

2.13 考虑如下形式的矩阵

$$A = I - \alpha B,$$

式中: α 为实数; $B \in \mathbb{R}^{n \times n}$ 为反对称矩阵, 即 $B^T = -B$. 证明:

(1) 对任意的非零向量 x , 有 $x^T A x / x^T x = 1$;

(2) 在 Arnoldi 分解中, 关于 V_k 的 Rayleigh 商 $H_k = V_k^T A V_k$ 的形式为

$$H_k = \begin{bmatrix} 1 & -\beta_1 & & \\ \beta_1 & 1 & \ddots & \\ & \ddots & \ddots & -\beta_{k-1} \\ & & \beta_{k-1} & 1 \end{bmatrix}.$$

2.14 证明 Kantorovich 不等式 (2.50).

第 3 章 线性方程组的矩阵分裂迭代法

在科学与工程计算中,很多问题往往最终归结为求解一个线性方程组问题,如结构分析、网络分析、大地测量、数据分析,以及用有限差分法或有限元法求解微分方程边值问题等.因此,研究求解大规模线性方程组快速、稳定的数值算法已成为当前科学与工程计算的核心问题之一.本章主要讨论求解线性方程组

$$Ax = b \quad (3.1)$$

的迭代方法,其中 A 是 n 阶可逆矩阵, b 是 n 维列向量.在实际应用中,矩阵 A 的阶数 n 一般很大,通常采用迭代法来求解.迭代法一般只涉及矩阵与向量的乘法运算.本章介绍迭代法的一些基本理论,以及基于矩阵分裂的几种经典迭代法,如 Jacobi 迭代法、Gauss-Seidel 迭代法以及松弛型迭代法、HSS 分裂迭代法等.当迭代矩阵的谱半径接近 1 时,迭代的收敛速度变得很缓慢,此时可用加速技巧对这些迭代法进行加速.

3.1 迭代法的一般理论

3.1.1 迭代法的定义与分类

迭代法本质上是一个递推公式,或者说是基于算子的重复利用.一般地,求解式 (3.1) 的迭代法定义如下.

定义 3.1 求解式 (3.1) 的迭代法就是寻找一个式 (3.1) 的近似解序列 $\{x^{(k)}\}_{k=0}^{\infty}$, 使得

$$\begin{cases} x^{(0)} = \phi_0(A, b), \\ x^{(k)} = \phi_k(x^{(0)}, x^{(1)}, \dots, x^{(k-1)}; A, b), \quad k = 1, 2, \dots, \end{cases} \quad (3.2)$$

式中: $\{x^{(k)}\}_{k=0}^{\infty}$ 称为迭代序列; $\{\phi_k\}_{k=0}^{\infty}$ 称为迭代算子序列; k 为迭代指标或迭代步数.

定义 3.2 若对某个整数 $s > 0$, $k \geq s$ 时, ϕ_k 与 k 无关,则称此迭代法为定常的,否则称为非定常的.若对任意的 k , ϕ_k 是 $x^{(0)}, x^{(1)}, \dots, x^{(k-1)}$ 的线性函数,则称迭代法是线性的,否则称为非线性的.

按照定义 3.2, 可以将迭代法分为线性定常迭代法、线性非定常迭代法、非线性定常迭代法、非线性非定常迭代法四类.

构造迭代法的一个基本思路是将式 (3.1) 变成如下形式的同解不动点方程组:

$$x = Bx + f, \quad (3.3)$$

然后任取一个初始点 $x^{(0)}$, 由迭代公式

$$x^{(k+1)} = Bx^{(k)} + f, \quad k = 0, 1, \dots \quad (3.4)$$

产生迭代序列 $\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}$, 其中 B 称为迭代矩阵. 当然, 若 B 与迭代指标 k 有关, 则迭代为非定常迭代, 否则为定常迭代. 可以构造各种可能的迭代矩阵 B , 但必须保证迭代序列收敛 (收敛性问题), 且收敛到方程组 (3.1) 的解 \mathbf{x}^* (相容性问题), 并由此确定出有效的迭代法. 一般地, 迭代序列的收敛过程是无限的, 而实际计算中只能且只需得到满足精度要求的某个近似解, 因此还要求适当地选取收敛准则, 这样, 一个迭代法才算完整.

3.1.2 收敛性与收敛速度

设 $\{\mathbf{x}^{(k)}\}$ 是 \mathbb{C}^n 中的某个向量序列, $\mathbf{x}^* \in \mathbb{C}^n$ 是某个向量, 有如下定义:

定义 3.3 若按某种范数有 $\lim_{k \rightarrow \infty} \|\mathbf{x}^{(k)} - \mathbf{x}^*\| = 0$, 则称向量序列收敛于 \mathbf{x}^* , 并记为 $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^*$.

事实上, $\lim_{k \rightarrow \infty} \|\mathbf{x}^{(k)} - \mathbf{x}^*\| = 0$ 当且仅当对 $i = 1, 2, \dots, n$, $\lim_{k \rightarrow \infty} x_i^{(k)} = x_i^*$.

记线性方程组 (3.1) 的精确解为 \mathbf{x}^* , 第 k 步迭代的误差向量记为 $\mathbf{e}^{(k)} = \mathbf{x}^* - \mathbf{x}^{(k)}$, 相应的残差向量记为 $\mathbf{r}^{(k)} = \mathbf{b} - A\mathbf{x}^{(k)}$.

定义 3.4 设 $\{\mathbf{x}^{(k)}\}$ 是某个迭代法产生的迭代序列, 若对任意的初始点 $\mathbf{x}^{(0)}$, 序列 $\{\mathbf{x}^{(k)}\}$ 均收敛到一个与 $\mathbf{x}^{(0)}$ 无关的极限, 则称此迭代法是收敛的. 若在精确运算下, 存在某个正整数 l , 使得 $\mathbf{x}^{(l)} = \mathbf{x}^*$, 则称此迭代法是精确收敛的或有限终止的.

下面给出收敛速度的概念. 注意到迭代法 (3.4) 的误差传播方程为

$$\mathbf{e}^{(k)} = B\mathbf{e}^{(k-1)} = \dots = B^k\mathbf{e}^{(0)}, \quad (3.5)$$

这里, $\mathbf{e}^{(0)} = \mathbf{x}^{(0)} - \mathbf{x}^*$ 是解的初始近似 $\mathbf{x}^{(0)}$ 与精确解的误差. 注意到迭代法 (3.4) 的误差向量 $\mathbf{e}^{(k)}$ 与右端项 \mathbf{f} 无关, 故其收敛性和收敛速度完全由迭代矩阵 B 确定.

引进误差向量后, 迭代的收敛问题就等价于误差向量序列收敛于 $\mathbf{0}$ 的问题.

欲使迭代法 (3.4) 对任意的初始向量 $\mathbf{x}^{(0)}$ 都收敛, 误差向量 $\mathbf{e}^{(k)}$ 应对任意的初始误差 $\mathbf{e}^{(0)}$ 都收敛于零向量. 于是迭代法 (3.4) 对任意的初始向量都收敛的充分必要条件是

$$\lim_{k \rightarrow \infty} B^k = O. \quad (3.6)$$

定理 3.1 迭代法 (3.4) 对任意的初始向量 $\mathbf{x}^{(0)}$ 都收敛的充分必要条件是 $\rho(B) < 1$, 这里 $\rho(B)$ 表示 B 的谱半径.

证明 必要性. 设对初始向量 $\mathbf{x}^{(0)}$, 迭代法 (3.4) 是收敛的, 那么式 (3.6) 成立. 由定理 1.17 (1), 对于任意的矩阵范数, 有

$$\rho(B) \leq \|B\|.$$

若 $\rho(B) < 1$ 不成立, 即 $\rho(B) \geq 1$, 则

$$\|B^k\| \geq \rho(B^k) = [\rho(B)]^k \geq 1,$$

这与式 (3.6) 矛盾.

充分性. 若 $\rho(\mathbf{B}) < 1$, 则存在一个正数 ε , 使得

$$\rho(\mathbf{B}) + 2\varepsilon < 1.$$

根据定理 1.17 (2), 存在一种矩阵范数 $\|\mathbf{B}\|_\varepsilon$, 使

$$\|\mathbf{B}\|_\varepsilon < \rho(\mathbf{B}) + \varepsilon < 1 - \varepsilon.$$

故得

$$\|\mathbf{B}^k\|_\varepsilon \leq \|\mathbf{B}\|_\varepsilon^k < (1 - \varepsilon)^k,$$

从而当 $k \rightarrow \infty$ 时, $\|\mathbf{B}^k\|_\varepsilon \rightarrow 0$, 即 $\mathbf{B}^k \rightarrow 0$. 证毕. \square

由此可见, 迭代是否收敛仅与迭代矩阵的谱半径有关, 即仅与方程组的系数矩阵和迭代格式的构造有关, 而与方程组的右端向量 \mathbf{b} 及初始向量 $\mathbf{x}^{(0)}$ 无关.

推论 3.1 若迭代矩阵 \mathbf{B} 的某种算子范数 $\|\mathbf{B}\|_p < 1$, 则迭代法 (3.4) 收敛.

设 \mathbf{B} 有 n 个线性无关的特征向量 $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n$ (可作为 n 维线性空间的一组基), 相应的特征值为 $\lambda_1, \lambda_2, \dots, \lambda_n$. 那么在这组基下, 初始误差向量 $\mathbf{e}^{(0)}$ 可以线性表示为

$$\mathbf{e}^{(0)} = \sum_{i=1}^n \alpha_i \mathbf{u}_i,$$

将其代入式 (3.5), 得

$$\mathbf{e}^{(k)} = \mathbf{B}^k \mathbf{e}^{(0)} = \sum_{i=1}^n \alpha_i \mathbf{B}^k \mathbf{u}_i = \sum_{i=1}^n \alpha_i \lambda_i^k \mathbf{u}_i. \quad (3.7)$$

故当 $\rho(\mathbf{B}) < 1$ 时, 其值越小, 迭代法收敛越快. 因此, 用 $\rho(\mathbf{B})$ 来刻画迭代法 (3.4) 的收敛速度是合适的.

另外, 对任何向量范数及相应的矩阵范数, 根据式 (3.5), 利用范数的性质, 有

$$\|\mathbf{e}^{(k)}\|_p \leq \|\mathbf{B}^k\|_p \|\mathbf{e}^{(0)}\|_p.$$

可见, $\|\mathbf{B}^k\|_p (< 1)$ 也给出了一种尺度, 它表示经 k 次迭代后误差向量范数减少的程度. 由此可定义平均收敛速度为

$$R_k(\mathbf{B}) = -\frac{1}{k} \log_{10} \|\mathbf{B}^k\|_p.$$

可以证明, 若 $\rho(\mathbf{B}) < 1$, 则有

$$\lim_{k \rightarrow \infty} (\|\mathbf{B}^k\|_p)^{\frac{1}{k}} = \rho(\mathbf{B}).$$

由于这一性质, 可定义渐近收敛速度为

$$R(\mathbf{B}) = \lim_{k \rightarrow \infty} R_k(\mathbf{B}) = -\log_{10} \rho(\mathbf{B}).$$

根据这个定义, $\rho(B)$ 越小, $R(B)$ 越大, 迭代法收敛越快. 值得注意的是, 尽管 $R_k(B)$ 依赖于所用的范数, 但 $R(B)$ 却与 k 无关. 因此, 通常将 $R(B)$ 称为收敛速度.

如果要使误差向量的范数减少一个因子 10^{-k_1} ($k_1 > 0$), 需迭代 k 次, 则可粗略求得 k 是满足不等式

$$k \geq k_1 / R(B)$$

的最小正整数, 即所要迭代次数与收敛速度成反比.

迭代法 (3.4) 是收敛的, 还可以给出近似解与准确解的误差估计.

定理 3.2 设迭代法 (3.4) 的迭代矩阵 B 满足 $\|B\| = q < 1$, 则迭代法是收敛的, 且有误差估计式

$$\|x^{(k)} - x^*\| \leq \frac{q^k}{1-q} \|x^{(0)} - x^{(1)}\|, \quad (3.8)$$

$$\|x^{(k)} - x^*\| \leq \frac{q}{1-q} \|x^{(k)} - x^{(k-1)}\|. \quad (3.9)$$

证明 由式 (3.5), 有

$$\|x^{(k)} - x^*\| = \|e^{(k)}\| \leq \|B^k\| \cdot \|e^{(0)}\| \leq q^k \|e^{(0)}\|.$$

注意到 $x^* = (I - B)^{-1}f$, 于是

$$\begin{aligned} \|e^{(0)}\| &= \|x^{(0)} - x^*\| = \|x^{(0)} - (I - B)^{-1}f\| \\ &= \|(I - B)^{-1}[(I - B)x^{(0)} - f]\| \\ &= \|(I - B)^{-1}(x^{(0)} - x^{(1)})\| \\ &\leq \|(I - B)^{-1}\| \cdot \|x^{(0)} - x^{(1)}\|. \end{aligned}$$

因为 $\|B\| < 1$, 根据定理 1.19, 有

$$\|(I - B)^{-1}\| \leq \frac{1}{1 - \|B\|} = \frac{1}{1 - q},$$

于是

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

下证式 (3.9). 由于

$$\begin{aligned} e^{(k)} &= x^{(k)} - x^* = (Bx^{(k-1)} + f) - (Bx^* + f) \\ &= Bx^{(k-1)} - Bx^* = Bx^{(k-1)} - B(I - B)^{-1}f \\ &= B(I - B)^{-1}[(I - B)x^{(k-1)} - f] \\ &= B(I - B)^{-1}(x^{(k-1)} - x^{(k)}), \end{aligned}$$

利用定理 1.19, 对上式两边取范数即得式 (3.9). 证毕. □

在理论上,可用式 (3.8) 估计近似解达到某一精度所需要的迭代次数 (但由于 q 不易计算,故计算实践中很少使用). 而式 (3.9) 则表明,只要 $\|B\|$ 不很接近于 1,即可用 $\{x^{(k)}\}$ 的相邻两项之差的范数 $\|x^{(k)} - x^{(k-1)}\|$ 来估计 $\|x^{(k)} - x^*\|$ 的大小,这为用 $\|x^{(k)} - x^{(k-1)}\|$ 作为算法的终止准则值提供了理论上的依据.

下面简要介绍一下矩阵分裂迭代法的基本知识. 设矩阵 A 非奇异且具有分裂

$$A = M - N, \quad (3.10)$$

其中 M 是非奇异的,则方程组 (3.1) 等价于

$$Mx = Nx + b. \quad (3.11)$$

构造迭代法

$$Mx^{(k+1)} = Nx^{(k)} + b, \quad k = 0, 1, \dots \quad (3.12)$$

求解方程组 (3.1), 这种方法统称为矩阵分裂迭代法,许多迭代法都可以写成这样的形式,如后面要介绍的 Richardson 迭代法、Jacobi 迭代法、Gauss-Seidel 迭代法、SOR 迭代法等. 在迭代法 (3.12) 中,由于每一次迭代需要解一个以 M 为系数矩阵的方程组,故一般要求非奇异矩阵 M 的形式比较简单,如对角矩阵、三对角矩阵、上(下)三角形矩阵等. 从后面的讨论知道,当 M 是 A 的一个很好的近似时,迭代解将很快收敛到方程组 (3.1) 的真解.

显然,迭代法 (3.12) 的迭代矩阵为 $B = M^{-1}N$, 并且它是线性定常迭代法. 根据定理 3.1 可知线性定常迭代法收敛的充要条件是

$$\rho(B) < 1. \quad (3.13)$$

3.1.3 相容性和敏感性分析

假设由一个迭代法所得的迭代序列是收敛的,问题是它是否收敛到待求解方程组的解. 对于线性定常迭代法,只要适当构造迭代矩阵,这不是个太大的问题. 但对非定常迭代法,就必须回答这问题. 需要指出,对于非定常迭代法,每一步迭代矩阵的谱半径 $\rho(B_k) < 1$ 既不是收敛的必要条件,也不是收敛的充分条件.

例 3.1 考虑方程组 $Ax = b$, 其中

$$A = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}, \quad b = \begin{bmatrix} 0 \\ 1 \end{bmatrix},$$

其精确解为

$$x^* = \begin{bmatrix} -1/3 \\ 2/3 \end{bmatrix}.$$

非定常迭代法

$$\begin{cases} x^{(k+1)} = B_1 x^{(k)} + f_1, & k \text{ 为偶数,} \\ x^{(k+1)} = B_2 x^{(k)} + f_2, & k \text{ 为奇数,} \end{cases} \quad (3.14)$$

式中:

$$B_1 = \begin{bmatrix} 0 & 0 \\ 2 & 0 \end{bmatrix}, \quad f_1 = \begin{bmatrix} -1/3 \\ 4/3 \end{bmatrix}, \quad B_2 = \begin{bmatrix} 0 & 2 \\ 0 & 0 \end{bmatrix}, \quad f_2 = \begin{bmatrix} -5/3 \\ 2/3 \end{bmatrix}.$$

显然, 每个迭代法均收敛, 因为 $\rho(B_1) = \rho(B_2) = 0$. 然而, 有

$$B_k = \begin{cases} B_1(B_2B_1)^{k/2} = \begin{bmatrix} 0 & 0 \\ 2^{k+1} & 0 \end{bmatrix}, & k \text{ 为偶数,} \\ (B_2B_1)^{(k+1)/2} = \begin{bmatrix} 2^{k+1} & 0 \\ 0 & 0 \end{bmatrix}, & k \text{ 为奇数,} \end{cases}$$

从而式 (3.13) 不成立.

设

$$B_1 = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}, \quad f_1 = \begin{bmatrix} 1/3 \\ -2/3 \end{bmatrix}, \quad B_2 = \begin{bmatrix} 1/4 & 0 \\ 0 & 1/4 \end{bmatrix}, \quad f_2 = \begin{bmatrix} -1/4 \\ 1/2 \end{bmatrix}.$$

显然前一方法不收敛. 然而, 有

$$B_2B_1 = \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \end{bmatrix},$$

$$B_k = \begin{cases} B_1(B_2B_1)^{k/2} = \begin{bmatrix} (1/2)^{(k-2)/2} & 0 \\ 0 & (1/2)^{(k-2)/2} \end{bmatrix}, & k \text{ 为偶数,} \\ (B_2B_1)^{(k+1)/2} = \begin{bmatrix} (1/2)^{(k+1)/2} & 0 \\ 0 & (1/2)^{(k+1)/2} \end{bmatrix}, & k \text{ 为奇数,} \end{cases}$$

从而 $\lim_{k \rightarrow \infty} B_k = O$, 且非定常迭代法 (3.14) 收敛.

下面考虑一种非定常迭代法, 其每一步使用公式

$$x^{(k+1)} = B_{k+1}x^{(k)} + f_{k+1}, \quad k \geq 0. \quad (3.15)$$

显然, $x^{(k)}$ 可表示为

$$x^{(k)} = \tilde{B}_k x^{(0)} + \tilde{f}_k, \quad k \geq 1. \quad (3.16)$$

其中

$$\tilde{B}_k = B_k B_{k-1} \cdots B_1, \quad (3.17)$$

且

$$\tilde{f}_k = f_k + B_k f_{k-1} + B_k B_{k-1} f_{k-2} + \cdots + B_k B_{k-1} \cdots B_2 f_1. \quad (3.18)$$

定义 3.5 (相容性) (1) 若 $x^{(k)} = x^*$, 且对所有的 $k' \geq k$ 有 $x^{(k')} = x^{(k)} = x^*$, 则称非定常迭代法 (3.16) 与式 (3.1) 是相容的.

(2) 若由非定常迭代法 (3.16) 定义的序列 $x^{(0)}, x^{(1)}, \dots$ 收敛, 且收敛到式 (3.1) 的解 x^* , 则称非定常迭代法 (3.16) 反相容于式 (3.1).

(3) 若非定常迭代法 (3.16) 与式 (3.1) 既是相容的又是反相容的, 则称其与式 (3.1) 是完全相容的.

定理 3.3 假设非定常迭代法 (3.16) 由式 (3.15) 而得, 若对每个 $s = 1, 2, \dots$, 线性定常迭代

$$x^{(k+1)} = B_s x^{(k)} + f_s, \quad k \geq 0 \quad (3.19)$$

是相容的, 则非定常迭代法 (3.16) 是相容的. 反之, 若非定常迭代法 (3.16) 是相容的, 则对每个 s , 线性定常迭代法 (3.19) 是相容的.

定理 3.4 若 A 是非奇异的且非定常迭代法 (3.16) 是相容的, 则对每个 k , 有

$$\tilde{f}_k = (I - \tilde{B}_k) A^{-1} b.$$

定理 3.5 非定常迭代法 (3.16) 收敛的充要条件是: $\{\tilde{f}_k\}$ 收敛, 且 $\lim_{k \rightarrow \infty} \tilde{B}_k = O$.

3.1.4 几种常见的矩阵分裂

设 $A \in \mathbb{R}^{n \times n} (\mathbb{C}^{n \times n})$ 非奇异, 若存在矩阵 M 和 N 满足 $A = M - N$, 且 M 是可逆的, 则称 M 和 N 为 A 的一个分裂.

定理 3.6 设 $A = M - N$ 是非奇异矩阵 $A \in \mathbb{C}^{n \times n}$ 的一个分裂, 则

$$M^{-1} N A^{-1} = A^{-1} N M^{-1}, \quad (3.20)$$

矩阵 $M^{-1} N$ 与 $A^{-1} N$ 可交换, 且 $N M^{-1}$ 与 $N A^{-1}$ 也可交换. 从而, 矩阵 $M^{-1} N$ 与 $A^{-1} N$ (或 $N M^{-1}$ 与 $N A^{-1}$) 具有相同的特征向量.

证明 直接验证式 (3.20), 得

$$\begin{aligned} M^{-1} N A^{-1} &= M^{-1} (M A^{-1} - I) = A^{-1} - M^{-1} \\ &= A^{-1} (M - A) M^{-1} = A^{-1} N M^{-1}. \end{aligned}$$

由式 (3.20) 立即可得矩阵 $M^{-1} N$ 与 $A^{-1} N$ 及 $N M^{-1}$ 与 $N A^{-1}$ 可交换.

设 x 为 $A^{-1} N$ 对应于特征值 λ 的特征向量, 则有

$$A^{-1} N x = \lambda x \implies N x = \lambda A x = \lambda (M - N) x.$$

由此可得

$$M^{-1} N x = \frac{\lambda}{1 + \lambda} x, \quad \lambda \neq -1.$$

上式表明 x 也是矩阵 $M^{-1} N$ 对应于特征值 $\frac{\lambda}{1 + \lambda}$ 的特征向量. 证毕. □

定理 3.7 设 $A = M - N$ 是非奇异矩阵 $A \in \mathbb{C}^{n \times n}$ 的一个分裂, 则有

$$\lambda(M^{-1}N) = \frac{\lambda(A^{-1}N)}{1 + \lambda(A^{-1}N)}. \quad (3.21)$$

因此, 若 $\lambda(A^{-1}N) \in \mathbb{R}$, 则相应的特征值 $\lambda(M^{-1}N) \in \mathbb{R}$, 反之亦然. 若 $\lambda(A^{-1}N) \in \mathbb{C}$, 则相应的特征值 $\lambda(M^{-1}N) \in \mathbb{C}$, 反之亦然.

下面给出文献中常见的一些矩阵分裂.

定义 3.6 设 $A = M - N$ 是非奇异矩阵 $A \in \mathbb{C}^{n \times n}$ 的一个分裂.

- (1) 若 $M^{-1} \geq O$ 且 $N \geq O$, 则称为正规分裂.
- (2) 若 $M^{-1} \geq O$, $M^{-1}N \geq O$ 且 $NM^{-1} \geq O$, 则称为非负分裂.
- (3) 若 $\rho(M^{-1}N) = \rho(NM^{-1}) < 1$, 则称为收敛的分裂.

定理 3.8 设 $A = M - N$ 是非奇异矩阵 $A \in \mathbb{C}^{n \times n}$ 的一个分裂, 如果 $M^{-1}N$ 和 $A^{-1}N$ 均非负, 则分裂是收敛的且有

$$\rho(M^{-1}N) = \frac{\rho(A^{-1}N)}{1 + \rho(A^{-1}N)}. \quad (3.22)$$

定义 3.7 设 $A = M - N$ 是非奇异矩阵 $A \in \mathbb{C}^{n \times n}$ 的一个分裂.

- (1) 若 $M^H + N$ 是 Hermite 正定的, 且 $N \geq O$, 则称为 P 正规分裂.
- (2) 若 $\langle M \rangle - |N|$ 是单调的 (即 $(\langle M \rangle - |N|)^{-1} \geq O$), 则称为 H 分裂.
- (3) 若 $\langle A \rangle = \langle M \rangle - |N|$, 则称为 H 相容分裂.
- (4) 若 M 是一个 M 矩阵且 $N \geq O$, 则称为 M 分裂.

下面给出使几种常见分裂收敛的一些条件及一些比较定理, 这些知识在证明相应的分裂迭代法的收敛性方面是很有用的.

定理 3.9 (1) 设 $A \in \mathbb{R}^{n \times n}$ 且 $A = M - N$ 为正规分裂, 如果 $A^{-1} \geq O$, 则

$$\rho(M^{-1}N) = \frac{\rho(A^{-1}N)}{1 + \rho(A^{-1}N)} < 1. \quad (3.23)$$

反之, 若 $\rho(M^{-1}N) < 1$, 则 $A^{-1} \geq O$.

(2) 设 $A \in \mathbb{C}^{n \times n}$ 是 Hermite 的且 $A = M - N$ 为 P 正规分裂, 则 $\rho(M^{-1}N) < 1$ 当且仅当 A 是 Hermite 正定的.

定理 3.9 说明, 在满足定理的条件下, 以 $M^{-1}N$ 为迭代矩阵的分裂迭代法是收敛的.

定理 3.10 设 $A^{-1} \geq O$ 且 $A = M_1 - N_1 = M_2 - N_2$ 为正规分裂, 如果下面条件之一成立:

- (1) $N_1 \leq N_2$.

- (2) $M_1^{-1} \geq M_2^{-1}$.
 (3) $M_2^{-1}N_2 \geq N_1M_1^{-1} \geq O$.
 (4) $A^{-1}N_2 \geq N_1A^{-1} \geq O$.

则

$$\rho(M_1^{-1}N_1) \leq \rho(M_2^{-1}N_2). \quad (3.24)$$

特别地, 若 $A^{-1} \geq O$ 且 $N_2 \geq N_1$ 而 $N_2 \neq N_1$ (或 $M_1^{-1} > M_2^{-1}$), 则

$$\rho(M_1^{-1}N_1) < \rho(M_2^{-1}N_2). \quad (3.25)$$

定理 3.10 说明, 在满足定理的条件下, 以 $A = M_1 - N_1$ 构造的分裂迭代法比以 $A = M_2 - N_2$ 构造的分裂迭代法收敛得更快.

定理 3.11 若 $A = M - N$ 为 A 的一个 H 分裂, 则 A 和 M 都是 H 矩阵, 且 $\rho(M^{-1}N) \leq \rho(|M|^{-1}|N|) < 1$. 若分裂是一个 H 相容分裂, 则它一定是 H 分裂, 从而是收敛的分裂.

定理 3.12 设 $A = M - N$ 为 A 的一个 M 分裂, 则

- (1) 若 A 是不可约的, 则存在一个正向量 $x > 0$ 使得 $(M^{-1}N)x = \rho(M^{-1}N)x$.
 (2) $\rho(M^{-1}N) < 1$ 当且仅当 $A = M - N$ 是一个非奇异 M 分裂.

3.2 几种经典迭代法

本节给出求解 $Ax = b$ 的几个经典迭代法及相应的收敛性定理.

3.2.1 Richardson 迭代法

在式 (3.10) 中取 $M = I$, $N = I - A$, 则相应于式 (3.12) 的迭代法为

$$x^{(k)} = (I - A)x^{(k-1)} + b = x^{(k-1)} + r^{(k-1)}, \quad (3.26)$$

式中: $r^{(k-1)} = b - Ax^{(k-1)}$ 为前一步的残差.

这就是著名的 Richardson 迭代法. 式 (3.26) 基本上不具有实用性, 因为它要求迭代矩阵 $I - A$ 的谱半径 $\rho(I - A) < 1$, 这等价于 $\rho(A) < 2$. 但它对于加速技术、Krylov 子空间方法及多重网格法都具有理论意义.

用 $-A$ 乘以式 (3.26) 两边并加上 b , 有

$$r^{(k)} = (I - A)r^{(k-1)} = \dots = (I - A)^k r^{(0)} = p_k(A)r^{(0)}, \quad (3.27)$$

式中: $p_k(A)$ 为 A 的 k 次多项式.

式 (3.27) 表明, 如果 $\|I - A\|$ 比 1 小得多, 则可期望其具有快速的收敛性.

注意到式 (3.27) 将残差表示成 A 的一个多项式与初始残差的乘积, 对于标准的 Richardson 迭代法, 此多项式非常简单, 即 $p_k(A) = (I - A)^k$. 以后将会看到, 许多流行的迭代法都具有类似的性质: 第 k 步的残差 $r^{(k)}$ 可以表示为 A 的某个 k 次多项式

$p_k(\mathbf{A})$ 与初始残差 $\mathbf{r}^{(0)}$ 的乘积, 并称此多项式为残差多项式. 这种性质是得到迭代法收敛界的一个有力工具.

还可以引进参数化的 Richardson 迭代法. 考虑 $\alpha > 0$ 并将 \mathbf{A} 分裂为

$$\mathbf{A} = \frac{1}{\alpha} \mathbf{I} - \left(\frac{1}{\alpha} \mathbf{I} - \mathbf{A} \right),$$

则导致一种定常的参数化 Richardson 迭代:

$$\mathbf{x}^{(k)} = \mathbf{x}^{(k-1)} + \alpha \mathbf{r}^{(k-1)}. \quad (3.28)$$

若在每一步迭代取不同的参数 α_k , 则得到非定常的参数化 Richardson 迭代:

$$\mathbf{x}^{(k)} = \mathbf{x}^{(k-1)} + \alpha_k \mathbf{r}^{(k-1)}. \quad (3.29)$$

上式中的 α_k 可有多种选取方式, 如通过极小化 $\|\mathbf{r}^{(k)}\|$ 来选取. 注意到

$$\mathbf{r}^{(k)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k-1)} - \alpha_k \mathbf{A}\mathbf{r}^{(k-1)} = (\mathbf{I} - \alpha_k \mathbf{A})\mathbf{r}^{(k-1)},$$

容易推得

$$\mathbf{r}^{(k)} = p_k(\mathbf{A})\mathbf{r}^{(0)}, \quad \text{其中 } p_k(\mathbf{A}) = \prod_{i=1}^k (\mathbf{I} - \alpha_i \mathbf{A}).$$

注 3.1 从式 (3.26)、式 (3.28) 和式 (3.29) 可知, Richardson 迭代法新的近似解 $\mathbf{x}^{(k)}$ 是前一步近似解 $\mathbf{x}^{(k-1)}$ 的一个校正, 校正量或者是前一步的残差 $\mathbf{r}^{(k-1)}$, 或者是沿残差方向前进某个步长 α_k (即 $\alpha_k \mathbf{r}^{(k-1)}$). 这一思想在现代变分迭代法中得到了很好的继承和发展. 最具代表性的是共轭梯度 (CG) 法和多重网格 (MG) 法.

由定理 3.1 不难得到定常的参数化 Richardson 迭代法 (3.28) 的收敛性定理.

定理 3.13 (1) 设 \mathbf{A} 为对称正定矩阵且其特征值 $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n > 0$, 则定常的参数化 Richardson 迭代法 (3.28) 收敛的充要条件是 $\alpha < \frac{2}{\lambda_1}$.

(2) 对于定常的参数化 Richardson 迭代法 (3.28), 参数 α 的最优值为 $\alpha_{\text{opt}} = \frac{2}{\lambda_1 + \lambda_n}$, 此时,

$$\rho(\mathbf{I} - \alpha_{\text{opt}} \mathbf{A}) = \frac{\kappa(\mathbf{A}) - 1}{\kappa(\mathbf{A}) + 1},$$

式中: $\kappa(\mathbf{A}) = \frac{\lambda_1}{\lambda_n}$.

3.2.2 Jacobi 迭代法

将系数矩阵 $\mathbf{A} = (a_{ij})$ 分裂为

$$\mathbf{A} = \mathbf{D} - \mathbf{L} - \mathbf{U}, \quad (3.30)$$

式中: $D = \text{diag}(a_{11}, a_{22}, \dots, a_{nn})$,

$$L = - \begin{bmatrix} 0 & & & & \\ a_{21} & 0 & & & \\ a_{31} & a_{32} & 0 & & \\ \vdots & \vdots & \ddots & \ddots & \\ a_{n1} & a_{n2} & \cdots & a_{n,n-1} & 0 \end{bmatrix}, \quad U = - \begin{bmatrix} 0 & a_{12} & a_{13} & \cdots & a_{1n} \\ & 0 & a_{23} & \cdots & a_{2n} \\ & & \ddots & \ddots & \vdots \\ & & & \ddots & a_{n-1,n} \\ & & & & 0 \end{bmatrix}.$$

在式 (3.10) 中, 取

$$M = D, \quad N = L + U, \quad (3.31)$$

则迭代法 (3.4) 中的迭代矩阵和右端项分别为

$$B_J = D^{-1}(L + U), \quad f_J = D^{-1}b.$$

相应的迭代法为

$$x^{(k+1)} = B_J x^{(k)} + f_J, \quad k = 0, 1, 2, \dots, \quad (3.32)$$

式 (3.32) 称为 Jacobi 迭代法. 其分量形式为

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left(b_i - \sum_{j \neq i} a_{ij} x_j^{(k)} \right), \quad i = 1, 2, \dots, n; \quad k = 0, 1, \dots. \quad (3.33)$$

该迭代法具有与 $x^{(k+1)}$ 中分量的计算次序无关, 因此容易实现并行计算等优点.

为了便于编程, 下面给出 Jacobi 迭代法的具体算法步骤如下:

算法 3.1 (Jacobi 迭代法)

- 步 1, 取初始点 $x^{(0)}$, 精度要求 ε , 最大迭代次数 N , 置 $k := 0$.
- 步 2, 由式 (3.32) 或式 (3.33) 计算 $x^{(k+1)}$.
- 步 3, 若 $\|b - Ax^{(k+1)}\| / \|b\| \leq \varepsilon$, 则停算, 输出 $x^{(k+1)}$ 作为方程组的近似解.
- 步 4, 置 $x^{(k)} := x^{(k+1)}$, $k := k + 1$, 转步 2.

根据算法 3.1, 可编制 MATLAB 程序如下:

```
%Jacobi迭代法程序-mjacobi.m
function [x,k,err,time]=mjacobi(A,b,x,tol,max_it)
if nargin<5, max_it=1000; end
if nargin<4, tol=1.e-5; end
if nargin<3, x=zeros(size(b)); end
tic; bnorm2 = norm(b);
r=b-A*x; %计算初始残差r0=b-Ax
err=norm(r)/bnorm2;
if (err<tol), return; end
D=diag(diag(A));
```

```

for k=1:max_it, % 迭代开始
    x=D\((D-A)*x+b);
    r=b-A*x; %计算残差r=b-Ax
    err=norm(r)/bnrm2;
    if(err<=tol), break; end
end
time=toc;

```

下面的定理给出了 Jacobi 迭代法的一个充分必要条件.

定理 3.14 设 $A = (a_{ij}) \in \mathbb{R}^{n \times n}$ 是对称矩阵, 且 $a_{ii} > 0$ ($i = 1, 2, \dots, n$), 则 Jacobi 迭代法收敛的充要条件是 A 和 $2D - A$ 都是正定矩阵.

证明 记 $D^{\frac{1}{2}} = \text{diag}(\sqrt{a_{11}}, \sqrt{a_{22}}, \dots, \sqrt{a_{nn}})$, $D^{-\frac{1}{2}} = (D^{\frac{1}{2}})^{-1}$, 则

$$B_J = D^{-1}(D - A) = D^{-\frac{1}{2}}(I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}})D^{\frac{1}{2}},$$

即 B_J 与 $I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ 相似, 从而它们有相同的特征值. 再由 A 的对称性可得

$$(I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}})^T = I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$$

即 $I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ 的特征值都是实数, 从而 B_J 的特征值都是实数.

必要性. 设 Jacobi 迭代法收敛, 即 $\rho(B_J) < 1$, 则 $\rho(I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}) < 1$. 设实数 λ 是实对称矩阵 $D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ 的任一特征值, 则有

$$0 = \det(\lambda I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}) = (-1)^n \det[(1 - \lambda)I - (I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}})],$$

即实数 $1 - \lambda$ 是实对称矩阵 $I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ 的特征值, 从而有

$$|1 - \lambda| < 1 \quad \text{或者} \quad 0 < \lambda < 2.$$

因此 $D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ 是正定矩阵. 由于 A 与 $D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ 合同, 所以 A 也是正定矩阵.

再设 $D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ 的全体特征值为 $\lambda_1, \lambda_2, \dots, \lambda_n$, 则存在正交矩阵 P , 使得

$$P^T(D^{-\frac{1}{2}}AD^{-\frac{1}{2}})P = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n).$$

于是有

$$P^T(2I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}})P = \text{diag}(2 - \lambda_1, 2 - \lambda_2, \dots, 2 - \lambda_n).$$

由于 $2 - \lambda_i > 0$ ($i = 1, 2, \dots, n$), 故 $2I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ 正定. 注意到 $2D - A = D^{\frac{1}{2}}(2I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}})D^{\frac{1}{2}}$ 合同于 $2I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$, 可得 $2D - A$ 正定.

充分性. 设 A 和 $2D - A$ 都是正定矩阵, 则有

$$A \text{ 正定} \Rightarrow D^{-\frac{1}{2}}AD^{-\frac{1}{2}} \text{ 正定} \Rightarrow D^{-\frac{1}{2}}AD^{-\frac{1}{2}} \text{ 的特征值大于零,}$$

$$2D - A \text{ 正定} \Rightarrow 2I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}} = D^{-\frac{1}{2}}(2D - A)D^{-\frac{1}{2}} \text{ 正定}$$

$$\Rightarrow 2I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}} \text{ 的特征值大于零.}$$

设 μ 是 B_J 的任一特征值, 从而也是 $I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ 的特征值, 且为实数, 则有

$$\begin{aligned} 0 &= \det(\mu I - B_J) = \det[\mu I - (I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}})] \\ &= \begin{cases} (-1)^n \det[(1 - \mu)I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}] \\ \det[(\mu + 1)I - (2I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}})] \end{cases}, \end{aligned}$$

即 $1 - \mu$ 是 $D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ 的特征值, $\mu + 1$ 是 $2I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ 的特征值, 所以

$$1 - \mu > 0, \mu + 1 > 0 \implies |\mu| < 1 \implies \rho(B_J) < 1.$$

由定理 3.1 知 Jacobi 迭代法收敛. 证毕. \square

下面的定理是 Jacobi 迭代法的另一个收敛性定理.

定理 3.15 若 A 为严格对角占优, 或不可约对角占优, 或 A 为 H 矩阵, 则 Jacobi 迭代法收敛.

证明 只证明 A 是严格对角占优或不可约对角占优的情形. 注意到 Jacobi 迭代法的迭代矩阵为

$$B_J = D^{-1}(D - A) = I - D^{-1}A.$$

只需证明 $\rho(B_J) < 1$. 事实上, 设 λ 为 B_J 的特征值, 对于行的情形, 当 A 为严格对角占优或不可约对角占优时, 矩阵 $I - B_J = D^{-1}A$ 也为严格对角占优或不可约对角占优. 因此, 当 $|\lambda| \geq 1$ 时, $I - \frac{1}{\lambda}B_J$ 也是严格对角占优或不可约对角占优的 (相当于 $I - B_J$ 的对角元不变, 非对角元模变小). 根据引理 1.1, 当 $|\lambda| \geq 1$ 时, 有

$$\det\left(I - \frac{1}{\lambda}B_J\right) \neq 0 \implies \det(\lambda I - B_J) \neq 0,$$

故 B_J 的特征值 λ 不满足 $|\lambda| \geq 1$, 即 $\rho(B_J) < 1$, 从而 Jacobi 迭代法收敛.

对于列的情形, 当 A 为严格对角占优或不可约对角占优时, 矩阵

$$I - DB_JD^{-1} = AD^{-1}$$

也为严格对角占优或不可约对角占优. 因此, 当 $|\lambda| \geq 1$ 时,

$$I - \frac{1}{\lambda}DB_JD^{-1}$$

也是严格对角占优或不可约对角占优的. 根据引理 1.1, 当 $|\lambda| \geq 1$ 时, 有

$$\det\left(I - \frac{1}{\lambda}DB_JD^{-1}\right) \neq 0,$$

从而

$$\det\left(I - \frac{1}{\lambda}B_J\right) \neq 0 \implies \det(\lambda I - B_J) \neq 0,$$

故 B_J 的特征值 λ 不满足 $|\lambda| \geq 1$, 即 $\rho(B_J) < 1$, 从而 Jacobi 迭代法收敛. 证毕. \square

Jacobi 迭代法的一种推广是将其用于分块矩阵

$$A = \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1m} \\ A_{21} & A_{22} & \cdots & A_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ A_{m1} & A_{m2} & \cdots & A_{mm} \end{bmatrix}, \quad \begin{aligned} A_{ir} &\in \mathbb{R}^{n_i \times n_r}, \quad 1 \leq i, r \leq m, \\ n_1 + n_2 + \cdots + n_m &= n, \end{aligned} \quad (3.34)$$

式中: 每个主对角块 $A_{ii} (i = 1, 2, \dots, m)$ 为非奇异的方阵. 将 $x^{(k+1)}$, $x^{(k)}$ 和 b 作相应的分块, 有

$$A_{ii}x_i^{(k+1)} = b_i - \sum_{r=1}^{i-1} A_{ir}x_r^{(k)} - \sum_{r=i+1}^m A_{ir}x_r^{(k)}, \quad i = 1, 2, \dots, m, \quad k = 0, 1, \dots \quad (3.35)$$

称上述迭代法为块 Jacobi (BJ) 迭代法, 其中用 x_i 表示第 i 个块分量以示区别. 它也有上面类似的收敛性定理.

另一种推广是在 Jacobi 迭代法中引入外推法技术, 即进行如下迭代

$$\begin{cases} D\tilde{x}^{(k+1)} = (L + U)x^{(k)} + b, \\ x^{(k+1)} = \gamma\tilde{x}^{(k+1)} + (1 - \gamma)x^{(k)}, \end{cases} \quad k = 0, 1, \dots, \quad (3.36)$$

式中: γ 为外推因子. 或等价地, 有

$$\frac{1}{\gamma}Dx^{(k+1)} = \frac{1-\gamma}{\gamma}Dx^{(k)} + (L + U)x^{(k)} + b, \quad (3.37)$$

该方法称为外推 Jacobi 方法, 其中, 若 $\gamma > 1$ 也称为 JOR 方法. 它相应于将 A 分裂为

$$A = M - N = \frac{1}{\gamma}D - \left[\frac{1-\gamma}{\gamma}D + (L + U) \right].$$

关于外推 Jacobi 迭代法, 有下面的收敛性定理.

定理 3.16 (1) 设 A 为对称正定的, 若 $\frac{2}{\gamma}D - A$ 是对称正定的, 则迭代法 (3.37) 对任意 $x^{(0)}$ 都收敛.

(2) 设 A 对称且 D 为正定矩阵, 则 $\frac{2}{\gamma}D - A$ 是对称正定的充要条件是

$$0 < \gamma < \frac{2}{1 - \lambda_{\min}(D^{-1}(D - A))}.$$

3.2.3 Gauss-Seidel (GS) 迭代法

从前面的讨论可以看到, Jacobi 迭代法的主要优点是方法简单, 在一定条件下具有实用性, 比如它可用作建立其他一些迭代过程的辅助方法. 然而, Jacobi 迭代法并不总是收敛的, 收敛时通常速度也很缓慢. 下面以 Jacobi 迭代法为基础, 考虑关于方程组

(3.1) 的其他迭代法. 仍假设 A 的对角元素不为零. 从串行计算的角度来看, Jacobi 迭代法很自然地先计算分量 x_1 的新迭代值

$$x_1^{(k+1)} = \frac{1}{a_{11}} \left(b_1 - \sum_{r=2}^n a_{1r} x_r^{(k)} \right). \quad (3.38)$$

得到 $x_1^{(k+1)}$ 后, 在后面 $x_i^{(k+1)}$ 的计算中就有理由用 $x_1^{(k+1)}$ 取代 $x_1^{(k)}$. 依次对 $x_2^{(k+1)}$, $x_3^{(k+1)}$ 等也用得到的新值立即取代旧值. Gauss-Seidel 原理就是“一旦获得新信息便立即利用”. 按此原理, 有

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left(b_i - \sum_{r=1}^{i-1} a_{ir} x_r^{(k+1)} - \sum_{r=i+1}^n a_{ir} x_r^{(k)} \right), \quad i = 1, 2, \dots, n; \quad k = 0, 1, \dots \quad (3.39)$$

此方法称为 Gauss-Seidel 迭代法 (GS 迭代法). 沿用前面的记号, 取

$$M = D - L, \quad N = U, \quad (3.40)$$

则式 (3.39) 可以写成

$$x^{(k+1)} = (D - L)^{-1} U x^{(k)} + (D - L)^{-1} b, \quad k = 0, 1, 2, \dots \quad (3.41)$$

因此, Gauss-Seidel 迭代法的迭代矩阵 $B_S = (D - L)^{-1} U$, 常向量 $f_S = (D - L)^{-1} b$. A 的对角元素非零的假设确保了 $(D - L)^{-1}$ 的存在性.

与 Jacobi 迭代法相比, Gauss-Seidel 迭代法使用了最新已经计算的分量. 因此, 一般情况下 Gauss-Seidel 迭代法要比 Jacobi 迭代法有效, 且在编程时 Gauss-Seidel 迭代法只要一个数组就够了, 将新计算出来的分量及时覆盖旧的分量.

为了便于计算机编程实现, 给出 Gauss-Seidel 迭代法的具体算法步骤如下.

算法 3.2 (GS 迭代法)

步 1, 输入矩阵 A , 右端向量 b , 初始点 $x^{(0)}$, 精度要求 ε , 最大迭代次数 N , 置 $k := 0$.

步 2, 由式 (3.41) 或式 (3.39) 计算 $x^{(k+1)}$.

步 3, 若 $\|b - Ax^{(k+1)}\| / \|b\| \leq \varepsilon$, 则停算, 输出 $x^{(k+1)}$ 作为方程组的近似解.

步 4, 置 $x^{(k)} := x^{(k+1)}$, $k := k + 1$, 转步 2.

根据算法 3.2, 编制 MATLAB 程序如下:

```
%GS迭代法程序-mseidel.m
function [x,k,err,time]=mseidel(A,b,x,tol,max_it)
if nargin<5, max_it=1000; end
if nargin<4, tol=1.e-5; end
if nargin<3, x=zeros(size(b)); end
tic; bnrm2=norm(b);
r=b-A*x; %计算初始残差r0=b-Ax
```

```

err=norm(r)/bnrm2;
if (err<tol), return; end
DL=tril(A); U=-triu(A,1);
for k=1:max_it,    % 迭代开始
    x=DL\ (U*x+b);
    r=b-A*x;    %计算残差r=b-Ax
    err=norm(r)/bnrm2;
    if (err<=tol), break; end
end
time=toc;

```

下面考虑 Gauss-Seidel 迭代法的收敛性.

定理 3.17 设 A 为严格对角占优, 或不可约对角占优, 或 A 是 H 矩阵, 则 Gauss-Seidel 迭代法收敛.

证明 只证明 A 为严格对角占优, 或不可约对角占优时的结论. 只需证 Gauss-Seidel 迭代法的迭代矩阵 $B_S = (D - L)^{-1}U$ 满足 $\rho(B_S) < 1$ 即可. 事实上, 设 μ 为 B_S 的特征值, 对于行的情形, 注意到 $(D - L) - U = A$, 若 A 为严格对角占优或不可约对角占优, 当 $|\mu| \geq 1$ 时,

$$(D - L) - \frac{1}{\mu}U$$

也是严格对角占优或不可约对角占优. 根据引理 1.1, 当 $|\mu| \geq 1$ 时, 有

$$\det \left[(D - L) - \frac{1}{\mu}U \right] \neq 0 \implies \det [\mu I - (D - L)^{-1}U] \neq 0,$$

故 B_S 特征值 μ 不满足 $|\mu| \geq 1$, 即 $\rho(B_S) < 1$, 从而 Gauss-Seidel 迭代法收敛.

对于列的情形, 注意到矩阵

$$(I - LD^{-1}) - UD^{-1} = (D - L - U)D^{-1} = AD^{-1},$$

因此, 若 A 为严格对角占优或不可约对角占优, 当 $|\mu| \geq 1$ 时, 有

$$(I - LD^{-1}) - \frac{1}{\mu}UD^{-1}$$

也为严格对角占优或不可约对角占优. 根据引理 1.1, 当 $|\mu| \geq 1$ 时, 有

$$\det \left[(I - LD^{-1}) - \frac{1}{\mu}UD^{-1} \right] = \det \left[(D - L) - \frac{1}{\mu}U \right] \det(D^{-1}) \neq 0,$$

故

$$\det \left[(D - L) - \frac{1}{\mu}U \right] \neq 0 \implies \det[\mu I - (D - L)^{-1}U] \neq 0,$$

即 B_S 特征值 μ 不满足 $|\mu| \geq 1$, 即 $\rho(B_S) < 1$, 从而 Gauss-Seidel 迭代法收敛. 证毕. \square

可见, 定理 3.15 中 Jacobi 迭代法收敛的充分条件也是 Gauss-Seidel 迭代法收敛的充分条件. 然而, 由于 Gauss-Seidel 迭代法利用了最新迭代信息, 在两种方法都收敛的情形下, Gauss-Seidel 迭代法往往比 Jacobi 迭代法收敛更快, 这是因为 $D-L$ 用到了 A 的更多的信息, 或 $(D-L)^{-1}A$ 比 $D^{-1}A$ 更接近单位矩阵.

定理 3.18 设 A 是 Hermite 正定矩阵, 则 Gauss-Seidel 迭代法收敛.

证明 设 λ 是 B_S 的任一特征值, 对应的特征向量为 y . 由 $B_S y = \lambda y$, 得

$$(D-L)^{-1}Uy = \lambda y \implies Uy = \lambda(D-L)y.$$

两端同时左乘 y^H , 得

$$y^H U y = \lambda(y^H D y - y^H L y).$$

因为 $A = D - L - U$ 是 Hermite 矩阵, 故 $U = L^H$. 记

$$y^H D y = p, \quad y^H L y = c + id \quad (c, d \text{ 为实数}),$$

则 $y^H U y = c - id$. 再由 A 正定, 得

$$D \text{ 正定} \implies p = y^H D y > 0,$$

$$y^H(D-L-U)y = y^H A y > 0 \implies p - 2c > 0.$$

于是有

$$\lambda = \frac{c - id}{p - (c + id)}, \quad |\lambda|^2 = \frac{c^2 + d^2}{(p - c)^2 + d^2}.$$

由于

$$c^2 - (p - c)^2 = -p(p - 2c) < 0,$$

所以 $|\lambda|^2 < 1$, 从而 $\rho(B_S) < 1$, 故 Gauss-Seidel 迭代法收敛. 证毕. \square

对具有如式 (3.34) 结构的分块系数矩阵及相应的向量划分, 有 Gauss-Seidel 迭代法的一种推广:

$$A_{ii}x_i^{(k+1)} = b_i - \sum_{r=1}^{i-1} A_{ir}x_r^{(k+1)} - \sum_{r=i+1}^m A_{ir}x_r^{(k)}, \quad i = 1, 2, \dots, m; \quad k = 0, 1, 2, \dots, \quad (3.42)$$

此即块 Gauss-Seidel (BGS) 迭代法.

另一种重要推广是双步法. 在第 $k+1$ 步, 作 Gauss-Seidel 迭代并将结果记为 $x^{(k+\frac{1}{2})}$:

$$x^{(k+\frac{1}{2})} = (D-L)^{-1}Ux^{(k)} + (D-L)^{-1}b, \quad k = 0, 1, 2, \dots, \quad (3.43)$$

然后各方程按相反顺序利用 $x^{(k+\frac{1}{2})}$ 进行计算:

$$x_i^{(k+1)} = \left(b_i - \sum_{r=1}^{i-1} a_{ir}x_r^{(k+\frac{1}{2})} - \sum_{r=i+1}^n a_{ir}x_r^{(k+1)} \right) / a_{ii}, \quad i = n, n-1, \dots, 1; \quad k = 0, 1, 2, \dots. \quad (3.44)$$

上式可写成矩阵形式

$$\mathbf{x}^{(k+1)} = (\mathbf{D} - \mathbf{U})^{-1} \mathbf{L} \mathbf{x}^{(k+\frac{1}{2})} + (\mathbf{D} - \mathbf{U})^{-1} \mathbf{b}, \quad k = 0, 1, 2, \dots \quad (3.45)$$

\mathbf{L} 与 \mathbf{U} 在式 (3.45) 中的作用正好与式 (3.43) 相反. 将式 (3.43) 代入式 (3.45), 得

$$\mathbf{x}^{(k+1)} = (\mathbf{D} - \mathbf{U})^{-1} \mathbf{L} (\mathbf{D} - \mathbf{L})^{-1} \mathbf{U} \mathbf{x}^{(k)} + \mathbf{f}, \quad k = 0, 1, 2, \dots, \quad (3.46)$$

式中:

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

式 (3.46) 称为对称 Gauss-Seidel (SGS) 迭代法.

3.3 松弛型迭代法

3.3.1 SOR 迭代法

如前面对 Jacobi 迭代法每个迭代步引入外推参数 (松弛因子) 那样, 可在 Gauss-Seidel 迭代法每个迭代步引入松弛因子, 从而得到所谓的 SOR 迭代法.

对 Gauss-Seidel 迭代法的一种重要改进如下:

$$x_i^{(k+1)} = (1 - \omega)x_i^{(k)} + \omega \tilde{x}_i^{(k+1)}, \quad i = 1, 2, \dots; \quad k = 0, 1, 2, \dots, \quad (3.47)$$

式中: $\omega > 0$ 为松弛因子, 是一个可以适当选取的参数, 用来加快收敛速度; $\tilde{x}_i^{(k+1)}$ 为式 (3.39) 的右端, 将其代入式 (3.47), 得

$$x_i^{(k+1)} = (1 - \omega)x_i^{(k)} + \frac{\omega}{a_{ii}} \left(b_i - \sum_{r=1}^{i-1} a_{ir} x_r^{(k+1)} - \sum_{r=i+1}^n a_{ir} x_r^{(k)} \right),$$

$$i = 1, 2, \dots, n; \quad k = 0, 1, \dots \quad (3.48)$$

也可以由矩阵分裂导出迭代法的矩阵向量形式. 取

$$\mathbf{M} = \frac{1}{\omega} \mathbf{D} - \mathbf{L}, \quad \mathbf{N} = \frac{1 - \omega}{\omega} \mathbf{D} + \mathbf{U}, \quad (3.49)$$

则有

$$\left(\frac{1}{\omega} \mathbf{D} - \mathbf{L} \right) \mathbf{x}^{(k+1)} = \left(\frac{1 - \omega}{\omega} \mathbf{D} + \mathbf{U} \right) \mathbf{x}^{(k)} + \mathbf{b}, \quad (3.50)$$

或等价地, 有

$$\mathbf{x}^{(k+1)} = (\mathbf{D} - \omega \mathbf{L})^{-1} [(1 - \omega) \mathbf{D} + \omega \mathbf{U}] \mathbf{x}^{(k)} + \omega (\mathbf{D} - \omega \mathbf{L})^{-1} \mathbf{b}, \quad (3.51)$$

这样定义的迭代法称为逐次超松弛迭代法 (Successive Over Relaxation, SOR 迭代法), 参数 ω 称为松弛因子. 适当选取 ω , 可使 SOR 迭代法的收敛速度优于 Gauss-Seidel 迭代法. 尽管 $0 < \omega < 1$ 时应该称为低松弛, 但为了方便, 对于任意的 $\omega \in (0, 2)$ 均使用

超松弛这一术语. 注意到, 在实际计算中, 松弛因子 ω 的选取通常是十分困难的. 没有一个通用的选取规则, 只能就某些具有特殊结构的矩阵讨论最佳松弛因子的计算.

SOR 迭代法的迭代矩阵和右端向量分别为

$$B_\omega = (D - \omega L)^{-1}[(1 - \omega)D + \omega U], \quad f_\omega = \omega(D - \omega L)^{-1}b. \quad (3.52)$$

在 A 的对角元素均非零的条件下 $(D - \omega L)^{-1}$ 是存在的. 显然, 当 $\omega = 1$ 时, SOR 迭代法退化为 Gauss-Seidel 迭代法.

为了便于计算机编程实现, 下面给出 SOR 迭代法的具体算法步骤.

算法 3.3 (SOR 迭代法)

步 1, 输入矩阵 A , 右端向量 b , 初始点 $x^{(0)}$, 精度要求 ε , 最大迭代次数 N , 置 $k := 0$.

步 2, 由式 (3.48) 或式 (3.51) 计算 $x^{(k+1)}$.

步 3, 若 $\|b - Ax^{(k+1)}\|/\|b\| \leq \varepsilon$, 则停算, 输出 $x^{(k+1)}$ 作为方程组的近似解.

步 4, 置 $x^{(k)} := x^{(k+1)}$, $k := k + 1$, 转步 2.

根据算法 3.3, 编制 MATLAB 程序如下:

```
%SOR迭代法程序-msor.m
function [x,k,err,time]=msor(A,b,w,x,tol,max_it)
if nargin<6, max_it=1000; end
if nargin<5, tol=1.e-5; end
if nargin<4, x=zeros(size(b)); end
tic; bnorm2=norm(b);
r=b-A*x; %计算初始残差r0=b-Ax
err=norm(r)/bnorm2;
if (err<tol), return; end
D=diag(diag(A)); L=-tril(A,-1); U=-triu(A,1);
for k=1:max_it % 迭代开始
    x=(D-w*L)\(((1-w)*D+w*U)*x+w*b);
    r=b-A*x; %计算残差r=b-Ax
    err= norm(r)/bnorm2;
    if (err<=tol), break; end
end
time=toc;
```

关于 SOR 迭代法的收敛性, 有如下的定理.

定理 3.19 对于任何参数 ω , SOR 迭代法的迭代矩阵 B_ω 满足

$$\rho(B_\omega) \geq |1 - \omega|. \quad (3.53)$$

证明 在

$$\begin{aligned} B_{\omega} &= (D - \omega L)^{-1} [(1 - \omega)D + \omega U] \\ &= (I - \omega D^{-1}L)^{-1} [(1 - \omega)I + \omega D^{-1}U] \end{aligned}$$

中, $D^{-1}L$ 和 $D^{-1}U$ 分别是严格下三角和严格上三角矩阵, 所以

$$\det(B_{\omega}) = \det(I - \omega D^{-1}L)^{-1} \det[(1 - \omega)I + \omega D^{-1}U] = (1 - \omega)^n.$$

再由 B_{ω} 的 n 个特征值之积等于 $\det(B_{\omega})$, 得

$$\rho(B_{\omega}) \geq (|\det(B_{\omega})|)^{\frac{1}{n}} = |1 - \omega|.$$

证毕. □

定理 3.20 设 ω 为实参数, 若 SOR 迭代法收敛, 则 $\omega \in (0, 2)$.

证明 根据定理 3.1, SOR 迭代法收敛时, $\rho(B_{\omega}) < 1$. 再由定理 3.19, 得 $|1 - \omega| < 1$, 即 $0 < \omega < 2$. 证毕. □

定理 3.20 表明, 如果松弛参数 $\omega \notin (0, 2)$, 则 SOR 迭代法不收敛. 但反之结论一般不成立, 即 $\omega \in (0, 2)$ 不能保证 SOR 迭代法收敛. 然而, 对 Hermite 正定矩阵 (对称正定矩阵), 这一条件是充分必要的.

定理 3.21 设 A 是 Hermite 正定矩阵, 则 SOR 迭代法收敛的充分必要条件是松弛因子 $\omega \in (0, 2)$.

证明 仅证明充分性即可. 设 λ 是 B_{ω} 的任一特征值, 对应的特征向量为 y . 由 $B_{\omega}y = \lambda y$, 得

$$[(1 - \omega)D + \omega U]y = \lambda(D - \omega L)y.$$

两端同时左乘 y^H , 得

$$(1 - \omega)y^H D y + \omega y^H U y = \lambda(y^H D y - \omega y^H L y).$$

因为 $A = D - L - U$ 是 Hermite 矩阵, 故 $U = L^H$. 记

$$y^H D y = p, \quad y^H L y = c + id \quad (c, d \text{ 为实数}),$$

则 $y^H U y = c - id$. 再由 A 正定, 得

$$\begin{aligned} D \text{ 正定} &\implies p = y^H D y > 0, \\ y^H (D - L - U)y &= y^H A y > 0 \implies p - 2c > 0. \end{aligned}$$

于是, 当 $0 < \omega < 2$ 时, $p - \omega(c + id) \neq 0$, 且有

$$\lambda = \frac{(1 - \omega)p + \omega(c - id)}{p - \omega(c + id)}, \quad |\lambda|^2 = \frac{[(1 - \omega)p + \omega c]^2 + \omega^2 d^2}{(p - \omega c)^2 + \omega^2 d^2}.$$

由于

$$[(1-\omega)p + \omega c]^2 - (p - \omega c)^2 = -p\omega(2-\omega)(p-2c) < 0,$$

所以 $|\lambda|^2 < 1$, 从而 $\rho(B_\omega) < 1$, 故 SOR 迭代法收敛. 证毕. \square

定理 3.22 设 A 严格对角占优或不可约对角占优, 则当实参数 ω 满足 $0 < \omega \leq 1$ 时, SOR 迭代法收敛.

证明 设 A 按行不可约对角占优, 且 λ 为 B_ω 的特征值. 如果 λ 满足 $|\lambda| \geq 1$, 则由 $0 < \omega \leq 1$, 得 $1 - \omega - \lambda \neq 0$. 再由

$$\lambda = (\lambda + \omega - 1) + (1 - \omega),$$

得

$$|\lambda + \omega - 1| \geq |\lambda| - (1 - \omega)$$

及

$$\left| \frac{\omega}{\lambda + \omega - 1} \right| \leq \left| \frac{\omega\lambda}{\lambda + \omega - 1} \right| \leq \frac{|\omega\lambda|}{|\lambda| - (1 - \omega)} \leq \frac{\omega|\lambda|}{|\lambda| - (1 - \omega)|\lambda|} = 1.$$

因为 $A = D - L - U$ 按行不可约对角占优, 从而

$$D - \frac{\omega\lambda}{\lambda + \omega - 1}L - \frac{\omega}{\lambda + \omega - 1}U$$

也按行不可约对角占优, 根据引理 1.1, 得

$$\det \left(D - \frac{\omega\lambda}{\lambda + \omega - 1}L - \frac{\omega}{\lambda + \omega - 1}U \right) \neq 0,$$

即

$$\det [\lambda(D - \omega L) - (1 - \omega)D - \omega U] \neq 0,$$

或

$$\det (\lambda I - (D - \omega L)^{-1}[(1 - \omega)D + \omega U]) \neq 0.$$

故 B_ω 的特征值 λ 不满足 $|\lambda| \geq 1$, 即 $\rho(B_\omega) < 1$, 从而 SOR 迭代法收敛.

当 A 为另外三种情形 (按列不可约对角占优、按行 (列) 严格对角占优) 时, 采用和上面类似的方法可以推得 SOR 迭代法收敛. 注意: 列的情形需要考虑矩阵 $I - LD^{-1} - UD^{-1} = AD^{-1}$. 证毕. \square

定理 3.23 设 A 为 H 矩阵, 若松弛因子 ω 满足

$$0 < \omega < \frac{2}{1 + \rho(B_J)},$$

式中: B_J 为 Jacobi 迭代法的迭代矩阵. 则 SOR 迭代法对任意的初始点 $x^{(0)}$ 均收敛.

SOR 迭代法也可推广到分块矩阵的情形. 对具有如式 (3.34) 结构的分块系数矩阵及相应的向量划分, 有 SOR 迭代法的一种推广:

$$\begin{aligned} \mathbf{x}_i^{(k+1)} &= (1-\omega)\mathbf{x}_i^{(k)} + \omega\mathbf{A}_{ii}^{-1}\left(\mathbf{b}_i - \sum_{r=1}^{i-1}\mathbf{A}_{ir}\mathbf{x}_r^{(k+1)} - \sum_{r=i+1}^n\mathbf{A}_{ir}\mathbf{x}_r^{(k)}\right), \\ i &= 1, 2, \dots, n; \quad k = 0, 1, \dots \end{aligned} \quad (3.54)$$

此即块 SOR (BSOR) 迭代法.

特别地, 对于 2×2 分块的经典鞍点方程组

$$\mathbf{A}\mathbf{u} := \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^T & \mathbf{O} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \begin{bmatrix} \mathbf{f} \\ \mathbf{g} \end{bmatrix} := \mathbf{b}, \quad (3.55)$$

式中: $\mathbf{A} \in \mathbb{R}^{m \times m}$ 为对称正定阵; $\mathbf{B} \in \mathbb{R}^{m \times n}$ 为列满秩阵.

注意: 此时系数矩阵 \mathbf{A} 非奇异, 且 Schur 补 $\mathbf{S} = \mathbf{B}^T\mathbf{A}^{-1}\mathbf{B}$ 对称正定. 给出分裂 $\mathbf{A} = \mathbf{D} - \mathbf{L} - \mathbf{U}$, 其中

$$\mathbf{D} = \begin{bmatrix} \mathbf{A} & \mathbf{O} \\ \mathbf{O} & \mathbf{S} \end{bmatrix}, \quad \mathbf{L} = \begin{bmatrix} \mathbf{O} & \mathbf{O} \\ -\mathbf{B}^T & \mathbf{O} \end{bmatrix}, \quad \mathbf{U} = \begin{bmatrix} \mathbf{O} & -\mathbf{B} \\ \mathbf{O} & \mathbf{S} \end{bmatrix},$$

由此得到迭代公式

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

即有

$$\begin{bmatrix} \mathbf{A} & \mathbf{O} \\ \omega\mathbf{B}^T & \mathbf{S} \end{bmatrix} \begin{bmatrix} \mathbf{x}^{(k+1)} \\ \mathbf{y}^{(k+1)} \end{bmatrix} = \begin{bmatrix} (1-\omega)\mathbf{A} & -\omega\mathbf{B} \\ \mathbf{O} & (1-\omega)\mathbf{S} \end{bmatrix} \begin{bmatrix} \mathbf{x}^{(k)} \\ \mathbf{y}^{(k)} \end{bmatrix} + \omega \begin{bmatrix} \mathbf{f} \\ \mathbf{g} \end{bmatrix}.$$

整理可得

$$\begin{cases} \mathbf{x}^{(k+1)} = (1-\omega)\mathbf{x}^{(k)} + \omega\mathbf{A}^{-1}(\mathbf{f} - \mathbf{B}\mathbf{y}^{(k)}), \\ \mathbf{y}^{(k+1)} = (1-\omega)\mathbf{y}^{(k)} + \omega\mathbf{S}^{-1}(\mathbf{g} - \mathbf{B}^T\mathbf{x}^{(k+1)}), \end{cases} \quad k = 0, 1, \dots$$

在实际计算中, 可选取对称正定矩阵 \mathbf{P} 和 \mathbf{Q} 为 \mathbf{A} 和 Schur 补 \mathbf{S} 的近似矩阵, 即得到所谓的非精确类算法.

3.3.2 SSOR 迭代法

SOR 迭代法的计算公式为

$$(\mathbf{D} - \omega\mathbf{L})\mathbf{x}^{(k+1)} = ((1-\omega)\mathbf{D} + \omega\mathbf{U})\mathbf{x}^{(k)} + \omega\mathbf{b}, \quad k = 0, 1, 2, \dots \quad (3.56)$$

它的计算是依赖顺序的, 即它是按从第 1 个到第 n 个分量依次计算的. 一个自然的想法是改变计算顺序, 即按从第 n 个到第 1 个分量逆序计算, 只需将式 (3.56) 中的 \mathbf{L} 和 \mathbf{U} 互换位置即得

$$(\mathbf{D} - \omega\mathbf{U})\mathbf{x}^{(k+1)} = ((1-\omega)\mathbf{D} + \omega\mathbf{L})\mathbf{x}^{(k)} + \omega\mathbf{b}, \quad k = 0, 1, 2, \dots \quad (3.57)$$

这种技巧导致一个对称方案, 称其为对称 SOR 迭代法 (Symmetric Successive Over Relaxation, SSOR 迭代法). 具体表示为

$$\begin{cases} (D - \omega L)x^{(k+\frac{1}{2})} = ((1 - \omega)D + \omega U)x^{(k)} + \omega b, \\ (D - \omega U)x^{(k+1)} = ((1 - \omega)D + \omega L)x^{(k+\frac{1}{2})} + \omega b, \end{cases} \quad k = 0, 1, 2, \dots \quad (3.58)$$

若消去 $x^{(k+\frac{1}{2})}$, 即得到形如式 (3.4) 的迭代法

$$x^{(k+1)} = B_{\text{SSOR}}x^{(k)} + f_{\text{SSOR}}, \quad (3.59)$$

式中:

$$\begin{aligned} B_{\text{SSOR}} &= (D - \omega U)^{-1}[(1 - \omega)D + \omega L](D - \omega L)^{-1}[(1 - \omega)D + \omega U], \\ f_{\text{SSOR}} &= \omega(2 - \omega)(D - \omega U)^{-1}D(D - \omega L)^{-1}b. \end{aligned} \quad (3.60)$$

从矩阵分裂的观点, SSOR 迭代法对应于分裂 $A = M - N$, 其中

$$M = \frac{1}{\omega(2 - \omega)}(D - \omega L)D^{-1}(D + \omega U), \quad (3.61)$$

$$N = \frac{1}{\omega(2 - \omega)}((1 - \omega)D + \omega L)D^{-1}((1 - \omega)D + \omega U). \quad (3.62)$$

此外, 若记 B_ω 是式 (3.52) 中定义的 SOR 迭代矩阵, \bar{B}_ω 是由 B_ω 中的 L 和 U 互换得到的. 则迭代矩阵 B_{SSOR} 可表示为

$$B_{\text{SSOR}} = M^{-1}N = \bar{B}_\omega B_\omega, \quad (3.63)$$

为了便于计算机编程实现, 下面给出 SSOR 迭代法的具体算法步骤.

算法 3.4 (SSOR 迭代法)

步 1, 输入矩阵 A , 右端向量 b , 初始点 $x^{(0)}$, 精度要求 ε , 最大迭代次数 N , 置 $k := 0$.

步 2, 由式 (3.58) 计算 $x^{(k+1)}$.

步 3, 若 $\|b - Ax^{(k+1)}\|/\|b\| \leq \varepsilon$, 则停算, 输出 $x^{(k+1)}$ 作为方程组的近似解.

步 4, 置 $x^{(k)} := x^{(k+1)}$, $k := k + 1$, 转步 2.

根据算法 3.4, 编制 MATLAB 程序如下:

```
%SSOR迭代法程序-mssor.m
[x,k,err,time]=mssor(A,b,w,x,tol,max_it)
if nargin<6, max_it=1000; end
if nargin<5, tol=1.e-6; end
if nargin<4, x=zeros(size(b)); end
tic; bnorm2=norm(b);
```

```

r=b-A*x; %计算初始残差
err=norm(r)/bnrm2;
if (err<tol), return; end
D=diag(diag(A)); L=-tril(A,-1); U=-triu(A,1);
for k=1:max_it %迭代开始
    x=(D-w*L)\(((1-w)*D+w*U)*x+w*b);
    x=(D-w*U)\(((1-w)*D+w*L)*x+w*b);
    r=b-A*x; %计算残差
    err=norm(r)/bnrm2;
    if (err<= tol), break; end
end
time=toc;

```

关于 SSOR 迭代法, 有下面的收敛性定理.

定理 3.24 SSOR 迭代法收敛的一个必要条件是 $|\omega - 1| < 1$. 对于 $\omega \in \mathbb{R}$, 条件变为 $\omega \in (0, 2)$.

定理 3.25 设 $A \in \mathbb{C}^{n \times n}$ 为具有正对角元的 Hermite 矩阵. 则对任意的 $\omega \in (0, 2)$, SSOR 迭代矩阵 B_{SSOR} 具有实非负特征值. 此外, 若 A 是对称正定的, 则 SSOR 迭代法收敛. 反之, 若 SSOR 迭代法收敛且 $\omega \in \mathbb{R}$, 则 $\omega \in (0, 2)$ 且 A 是对称正定的.

注 3.2 SSOR 迭代法实质上就是将 L 和 U 等同看待连续两次使用 SOR 迭代法. 这样做好处有两个:

- (1) 某些特殊问题, 用 SOR 迭代法不收敛, 但依然可构造出收敛的 SSOR 迭代法.
- (2) 一般来说, SOR 迭代法的渐近收敛速度对松弛因子 ω 的选择是非常敏感的, 而 SSOR 却不敏感.

SSOR 迭代法的一种推广是它的两个半步迭代格式中使用不同的松弛因子:

$$\begin{cases} (D - \omega_1 L)x^{(k+\frac{1}{2})} = ((1 - \omega_1)D + \omega_1 U)x^{(k)} + \omega_1 b, \\ (D - \omega_2 U)x^{(k+1)} = ((1 - \omega_2)D + \omega_2 L)x^{(k+\frac{1}{2})} + \omega_2 b, \end{cases} \quad k = 0, 1, 2, \dots \quad (3.64)$$

若消去 $x^{(k+\frac{1}{2})}$, 得

$$x^{(k+1)} = B_{\omega_1, \omega_2} x^{(k)} + f_{\omega_1, \omega_2}, \quad (3.65)$$

式中:

$$\begin{aligned} B_{\omega_1, \omega_2} &= (D - \omega_2 U)^{-1}[(1 - \omega_2)D + \omega_2 L](D - \omega_1 L)^{-1}[(1 - \omega_1)D + \omega_1 D], \\ f_{\omega_1, \omega_2} &= (\omega_1 + \omega_2 - \omega_1 \omega_2)(D - \omega_2 U)^{-1}D(D - \omega_1 L)^{-1}b. \end{aligned} \quad (3.66)$$

这种推广得到的迭代法通常称为不对称 SOR 迭代法, 简记为 USSOR 迭代法.

进一步, SSOR 迭代法也可推广到分块矩阵的情形. 对具有如式 (3.34) 结构的分块系数矩阵及相应的向量划分, 有 SSOR 迭代法的另一种推广:

$$\begin{aligned} \mathbf{x}_i^{(k+\frac{1}{2})} &= (1-\omega)\mathbf{x}_i^{(k)} + \omega\mathbf{A}_{ii}^{-1}\left(\mathbf{b}_i - \sum_{r=1}^{i-1}\mathbf{A}_{ir}\mathbf{x}_r^{(k+1)} - \sum_{r=i+1}^n\mathbf{A}_{ir}\mathbf{x}_r^{(k)}\right), \quad i=1,2,\cdots,n, \\ \mathbf{x}_i^{(k+1)} &= (1-\omega)\mathbf{x}_i^{(k+\frac{1}{2})} + \omega\mathbf{A}_{ii}^{-1}\left(\mathbf{b}_i - \sum_{r=1}^{i-1}\mathbf{A}_{ir}\mathbf{x}_r^{(k+\frac{1}{2})} - \sum_{r=i+1}^n\mathbf{A}_{ir}\mathbf{x}_r^{(k+1)}\right), \\ i &= n, n-1, \cdots, 1; \quad k=0, 1, 2, \cdots. \end{aligned} \quad (3.67)$$

此即块 SSOR (BSSOR) 迭代法.

例 3.2 用不同的迭代法求解 n 元线性方程组 $\mathbf{Ax} = \mathbf{b}$, 其中

$$\mathbf{A} = \begin{bmatrix} 4 & -1 & & & \\ -1 & 4 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 4 & -1 \\ & & & -1 & 4 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 3 \\ 2 \\ \vdots \\ 2 \\ 3 \end{bmatrix}. \quad (3.68)$$

方程的精确解为 $\mathbf{x}^* = (1, 1, \cdots, 1)^T$. 取 $n = 2^{12} - 1 = 4095$, 初始向量为零向量, 容许误差为 10^{-10} .

解 用 Jacobi, GS, SOR 和 SSOR 四种迭代法进行测试, 取松弛因子 $\omega = 1.1$, 得到计算结果如表 3.1 所示.

表 3.1 Jacobi, GS, SOR 和 SSOR 四种迭代法的数值比较

迭代法	迭代次数(k)	CPU 时间	相对残差 ($\ \mathbf{r}^{(k)}\ _2/\ \mathbf{b}\ _2$)
Jacobi	34	0.0043	5.8104e-11
GS	21	0.0030	9.5383e-11
SOR	17	0.0049	3.4644e-11
SSOR	9	0.0051	8.0601e-12

3.3.3 AOR 迭代法

有一种技术可以对收敛的迭代格式进行“加速”, 或者使得不收敛的迭代格式变得收敛. 这种“加速”技术通常是通过引进一个“加速”或“松弛”参数 $\gamma \in \mathbb{C} \setminus \{0\}$ 来实现的. 基于矩阵分裂 $\mathbf{A} = \mathbf{M} - \mathbf{N}$, 考虑一个新的矩阵 $\mathbf{M}_\gamma = \frac{1}{\gamma}\mathbf{M}$. 设原迭代法的迭代格式为 $\mathbf{x}^{(k+1)} = \mathbf{B}\mathbf{x}^{(k)} + \mathbf{f}$, 则新迭代法的迭代格式为

$$\mathbf{x}^{(k+1)} = \mathbf{B}_\gamma\mathbf{x}^{(k)} + \mathbf{f}_\gamma, \quad \mathbf{B}_\gamma = (1-\gamma)\mathbf{I} + \gamma\mathbf{B}, \quad \mathbf{f}_\gamma = \gamma\mathbf{f}.$$

参数 $\gamma \in \mathbb{C} \setminus \{0\}$ 称为外推参数, 且相应的格式称为原格式的外推. 最优外推参数的确定一般需要较强的假设条件并涉及原迭代矩阵 B 的谱 $\sigma(B)$ 的有关信息.

利用外推法思想, Hadjidimas (哈吉迪马斯) 于 1978 年引进一种具有两个松弛参数 r 和 ω 的 SOR 类迭代法, 称为加速超松弛方法 (简记为 AOR 方法), 其定义如下:

$$x^{(k+1)} = B_{r,\omega} x^{(k)} + f_{r,\omega}, \quad (3.69)$$

式中:

$$B_{r,\omega} = (D - rL)^{-1} [(1 - \omega)D + (\omega - r)L + \omega U], \quad f_{r,\omega} = \omega(D - rL)^{-1}b. \quad (3.70)$$

容易证明: 当 $r = 0$ 时, AOR 方法就是具有外推参数 ω 的外推 Jacobi 方法; 当 $r \neq 0$ 时, AOR 方法就是具有外推参数 $s = \omega/r$ 的外推 SOR 方法, 其中 r 是原 SOR 方法的松弛因子. 此外, 显然前面的 Jacobi, Gauss-Seidel, SOR 方法等均可认为是 AOR 方法的特殊情形:

- (1) 当 $r = 0, \omega = 1$ 时, AOR 方法即为 Jacobi 方法.
- (2) 当 $r = \omega = 1$ 时, AOR 方法即为 Gauss-Seidel 方法.
- (3) 当 $r = \omega \neq 0$ 时, AOR 方法即为 SOR 方法.

关于 AOR 方法, 有下面的收敛性定理.

定理 3.26 设 $A \in \mathbb{C}^{n \times n}$ 为 Hermite 矩阵, 且 $A = D - L - L^H$, D 是 Hermite 正定的, $\det(D - rL) \neq 0, \forall \omega \in (0, 2)$ 且 $r \in (\omega + (2 - \omega)/\mu_{\min}, \omega + (2 - \omega)/\mu_{\max})$, 其中 $\mu_{\min} < 0 < \mu_{\max}$ 为矩阵 $D^{-1}(L + L^H)$ 的最小和最大特征值. 则 $\rho(B_{r,\omega}) < 1$ 当且仅当 A 是正定的.

类似于 SSOR 方法, 可以推广 AOR 方法得到对称 AOR 方法 (简记为 SAOR 方法). 具体形式为

$$\begin{cases} (D - rL)x^{(k+\frac{1}{2})} = [(1 - \omega)D + (\omega - r)L + \omega U]x^{(k)} + \omega b, \\ (D - rU)x^{(k+1)} = [(1 - \omega)D + (\omega - r)U + \omega L]x^{(k+\frac{1}{2})} + \omega b, \end{cases} \quad k = 0, 1, 2, \dots$$

上式可以化简为

$$\begin{cases} (D - rL)x^{(k+\frac{1}{2})} = (D - rL - \omega A)x^{(k)} + \omega b, \\ (D - rU)x^{(k+1)} = (D - rU - \omega A)x^{(k+\frac{1}{2})} + \omega b, \end{cases} \quad k = 0, 1, 2, \dots \quad (3.71)$$

消去 $x^{(k+\frac{1}{2})}$, 得

$$x^{(k+1)} = B_{\text{SAOR}} x^{(k)} + f_{\text{SAOR}}, \quad (3.72)$$

式中:

$$\begin{aligned} B_{\text{SAOR}} &= (D - rL)^{-1}(D - rU - \omega A)(D - rU)^{-1}(D - rL - \omega A), \\ f_{\text{SAOR}} &= \omega(D - rU)^{-1}(2D - rL - \omega A)(D - rL)^{-1}b. \end{aligned}$$

同样可以将 AOR 方法推广到非对称 AOR 方法 (UAOR 方法) 和分块 AOR 方法 (BAOR 方法) 的情形, 这里不再详述.

3.4 HSS 迭代法

本节内容主要取材于文献 [21, 22]. HSS 迭代法是中国学者白中治等提出的用于求解大型稀疏非 Hermite 正定方程组的一种有效的数值方法. 它是基于系数矩阵的 Hermite 和反 Hermite 分裂 (Hermitian/skew-Hermitian splitting, HSS), 其中包括 HSS 迭代和非精确 HSS (IHSS) 迭代, IHSS 在 HSS 外迭代的每一步用某种 Krylov 子空间方法 (见第 4 章) 求解子问题作为其内迭代过程. 理论分析显示 HSS 方法无条件地收敛到线性方程组的唯一解. 目前, 基于 HSS 的各种迭代方法及与不同预处理方式的结合成为一个研究热点.

3.4.1 HSS 和 IHSS 方法

许多科学计算问题需要求解 (复) 线性方程组

$$Ax = b, \quad (3.73)$$

式中: $A \in \mathbb{C}^{n \times n}$ 为非奇异的大型稀疏非 Hermite 正定矩阵; $x, b \in \mathbb{C}^n$.

将矩阵 A 进行 Hermite-反 Hermite 分裂 (HS) 分裂

$$A = H + S, \quad (3.74)$$

式中:

$$H = \frac{1}{2}(A + A^H), \quad S = \frac{1}{2}(A - A^H). \quad (3.75)$$

下面将研究基于此特殊矩阵分裂的求解线性方程组 (3.73) 的有效迭代方法.

基于 HS 分裂式 (3.74), 文献 [21] 提出了求解线性方程组 (3.73) 的 Hermite-反 Hermite 分裂迭代法 (简称 HSS 迭代法), 其形式如下:

算法 3.5 (HSS 迭代法)

步 1, 给定初始向量 $x^{(0)}$, 选取参数 $\alpha > 0$, 容许误差 $\varepsilon > 0$, 置 $k := 0$.

步 2, 计算

$$\begin{cases} (\alpha I + H)x^{(k+\frac{1}{2})} = (\alpha I - S)x^{(k)} + b, \\ (\alpha I + S)x^{(k+1)} = (\alpha I - H)x^{(k+\frac{1}{2})} + b, \end{cases}$$

步 3, 计算残差 $r^{(k+1)} = b - Ax^{(k+1)}$. 若 $\|r^{(k+1)}\|_2 \leq \varepsilon$, 停算. 否则, 置 $k := k + 1$, 转步 2.

明显地, HSS 迭代法的每一步都在矩阵 A 的 Hermite 部分和反 Hermite 部分之间进行交替, 这类似于求解偏微分方程的交替方向隐式 (ADI) 迭代法.

注意到可以调换上述 HSS 迭代法中矩阵 H 和 S 的角色, 即先解关于 $\alpha I + S$ 的方程组, 然后再解关于 $\alpha I + H$ 的方程组. 每个 HSS 迭代的两个半步需要精确求解具有 n 阶系数矩阵 $\alpha I + H$ 和 $\alpha I + S$ 的方程组. 然而, 这在实际实现中是耗时和不切实际的. 为了改进 HSS 迭代法的计算效率, 例如, 在 HSS 迭代法的每一步利用共轭梯度法 (CG 法) 求解关于 $\alpha I + H$ 的方程组 (因为 $\alpha I + H$ 是 Hermite 正定的), 并用某种 Krylov 子空间方法 (如 GMRES 方法, 见第 4 章) 解关于 $\alpha I + S$ 的方程组, 直到预

先设定的精度. 这导致一种非精确的 Hermite-反 Hermite 分裂迭代法, 简称 IHSS 迭代法. 内迭代的阈值 (或内迭代步数) 可随外迭代方案的不同而变化. 因此, IHSS 迭代实际上是求解式 (3.73) 的一种非定常迭代法. 对给定正常数 α , IHSS 迭代法形式如下.

算法 3.6 (IHSS 迭代法)

步 1, 给定初始向量 $\mathbf{x}^{(0)}$, 选取参数 $\alpha > 0$, 容许误差 $\varepsilon > 0$, 置 $k := 0$.

步 2, 以 $\mathbf{x}^{(k)}$ 为初值, 利用一种内迭代 (如 CG 法) 近似求解 $\mathbf{x}^{(k+\frac{1}{2})}$:

$$(\alpha I + H)\mathbf{x}^{(k+\frac{1}{2})} \approx (\alpha I - S)\mathbf{x}^{(k)} + \mathbf{b}.$$

步 3, 以 $\mathbf{x}^{(k+\frac{1}{2})}$ 为初值, 利用一种内迭代 (如某种 krylov 子空间方法) 近似求解 $\mathbf{x}^{(k+1)}$:

$$(\alpha I + S)\mathbf{x}^{(k+1)} \approx (\alpha I - H)\mathbf{x}^{(k+\frac{1}{2})} + \mathbf{b}.$$

步 4, 计算残差 $\mathbf{r}^{(k+1)} = \mathbf{b} - A\mathbf{x}^{(k+1)}$. 若 $\|\mathbf{r}^{(k+1)}\|_2 \leq \varepsilon$, 停算. 否则, 置 $k := k+1$, 转步 2.

为了便于数值实现和理论分析, 可将上面的 IHSS 迭代法重写为如下等价形式: 其中 $\|\cdot\|$ 是某种向量范数, 并记原方程组的残差为 $\mathbf{r}^{(k)} = \mathbf{b} - A\mathbf{x}^{(k)}$, $\mathbf{r}^{(k+\frac{1}{2})} = \mathbf{b} - A\mathbf{x}^{(k+\frac{1}{2})}$, 内迭代的残差分别记为

$$\mathbf{p}^{(k)} = \mathbf{r}^{(k)} - (\alpha I + H)\mathbf{z}^{(k)}; \quad (3.76)$$

$$\mathbf{q}^{(k+\frac{1}{2})} = \mathbf{r}^{(k+\frac{1}{2})} - (\alpha I + S)\mathbf{z}^{(k+\frac{1}{2})}. \quad (3.77)$$

则可以改写算法 3.6 为如下形式.

算法 3.7 (IHSS 迭代法)

步 1, 给定初始向量 $\mathbf{x}^{(0)}$, 选取参数 $\alpha > 0$, 容许误差 $\varepsilon > 0$, 置 $k := 0$.

步 2, 以 $\mathbf{x}^{(k)}$ 为初值, 迭代求解 $(\alpha I + H)\mathbf{z}^{(k)} = \mathbf{r}^{(k)}$, 直到 $\|\mathbf{p}^{(k)}\| \leq \varepsilon_k \|\mathbf{r}^{(k)}\|$.

步 3, 计算 $\mathbf{x}^{(k+\frac{1}{2})} = \mathbf{x}^{(k)} + \mathbf{z}^{(k)}$.

步 4, 以 $\mathbf{x}^{(k+\frac{1}{2})}$ 为初值, 迭代求解 $(\alpha I + S)\mathbf{z}^{(k+\frac{1}{2})} = \mathbf{r}^{(k+\frac{1}{2})}$, 直到 $\|\mathbf{q}^{(k+\frac{1}{2})}\| \leq \eta_k \|\mathbf{r}^{(k+\frac{1}{2})}\|$.

步 5, 计算 $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k+\frac{1}{2})} + \mathbf{z}^{(k+\frac{1}{2})}$.

步 6, 计算残差 $\mathbf{r}^{(k+1)} = \mathbf{b} - A\mathbf{x}^{(k+1)}$. 若 $\|\mathbf{r}^{(k+1)}\|_2 \leq \varepsilon$, 停算. 否则, 置 $k := k+1$, 转步 2.

事实上, 这里构造了求解非 Hermite 正定方程组的一个一般的迭代法框架. 在此迭代框架下有一些组合. 可精确或非精确地求解 Hermite 部分, 也可精确或非精确地求解反 Hermite 部分. 最佳选择依赖于 Hermite 和反 Hermite 矩阵的结构. 所得的相应精确或非精确 Hermite-反 Hermite 分裂迭代的收敛性理论可类似于下面的分析来建立, 只要做少许的修改即可.

下面考虑 HSS 和 IHSS 迭代法的收敛性和收敛速度方面的一些理论结果. 首先注意到 HSS 迭代法可推广到两步分裂迭代框架, 下面的引理先给出两步分裂迭代的一个一般性收敛条件.

引理 3.1 设 $A \in \mathbb{C}^{n \times n}$, $A = M_i - N_i$ ($i = 1, 2$) 为矩阵 A 的两个分裂. $x^{(0)}$ 为给定的初始向量. 如果 $\{x^{(k)}\}$ 是由下面两步分裂迭代法所定义的序列

$$\begin{cases} M_1 x^{(k+\frac{1}{2})} = N_1 x^{(k)} + b, \\ M_2 x^{(k+1)} = N_2 x^{(k+\frac{1}{2})} + b, \end{cases} \quad k = 0, 1, 2, \dots$$

则

$$x^{(k+1)} = M_2^{-1} N_2 M_1^{-1} N_1 x^{(k)} + M_2^{-1} (M_1 + N_2) M_1^{-1} b, \quad k = 0, 1, 2, \dots$$

进一步, 如果迭代矩阵 $M_2^{-1} N_2 M_1^{-1} N_1$ 的谱半径 $\rho(M_2^{-1} N_2 M_1^{-1} N_1)$ 小于 1, 则迭代序列 $\{x^{(k)}\}$ 对所有初始向量 $x^{(0)} \in \mathbb{C}^n$ 都收敛于线性方程组 (3.73) 的唯一解 $x^* \in \mathbb{C}^n$.

对于 HSS 迭代法, 应用引理 3.1 及矩阵谱的相似不变性易得如下定理.

定理 3.27 设 $A \in \mathbb{C}^{n \times n}$ 为正定矩阵, $H = \frac{1}{2}(A + A^H)$ 和 $S = \frac{1}{2}(A - A^H)$ 为其 Hermite 和反 Hermite 部分, 且 α 是一个正常数. 则 HSS 迭代矩阵 $M(\alpha)$ 为

$$M(\alpha) = (\alpha I + S)^{-1}(\alpha I - H)(\alpha I + H)^{-1}(\alpha I - S), \quad (3.78)$$

且其谱半径 $\rho(M(\alpha))$ 有上界

$$\sigma(\alpha) \equiv \max_{\lambda_i \in \lambda(H)} \left| \frac{\alpha - \lambda_i}{\alpha + \lambda_i} \right|,$$

式中: $\lambda(H)$ 为矩阵 H 的谱集, 因此成立

$$\rho(M(\alpha)) \leq \sigma(\alpha) < 1, \quad \forall \alpha > 0,$$

即 HSS 迭代收敛到线性方程组 (3.73) 的唯一解 $x^* \in \mathbb{C}^n$.

证明 在引理 3.1 中, 记

$$M_1 = \alpha I + H, N_1 = \alpha I - S, M_2 = \alpha I + S, N_2 = \alpha I - H,$$

注意到对任意的 $\alpha > 0$, $\alpha I + H$ 和 $\alpha I + S$ 是非奇异的, 则由引理 3.1 立即可得式 (3.78).

进一步, 利用矩阵谱的相似不变性, 得

$$\begin{aligned} \rho(M(\alpha)) &= \rho((\alpha I + S)^{-1}(\alpha I - H)(\alpha I + H)^{-1}(\alpha I - S)) \\ &= \rho((\alpha I - H)(\alpha I + H)^{-1}(\alpha I - S)(\alpha I + S)^{-1}) \\ &\leq \|(\alpha I - H)(\alpha I + H)^{-1}\|_2 \cdot \|(\alpha I - S)(\alpha I + S)^{-1}\|_2. \end{aligned}$$

记 $Q(\alpha) = (\alpha I - S)(\alpha I + S)^{-1}$, 注意到 $S^H = -S$, 有

$$\begin{aligned} Q(\alpha)^H Q(\alpha) &= (\alpha I - S)^{-1}(\alpha I + S)(\alpha I - S)(\alpha I + S)^{-1} \\ &= (\alpha I - S)^{-1}(\alpha I - S)(\alpha I + S)(\alpha I + S)^{-1} = I, \end{aligned}$$

即 $Q(\alpha)$ 是酉矩阵, 故有 $\|Q(\alpha)\|_2 = 1$. 从而, 有

$$\rho(M(\alpha)) \leq \|(\alpha I - H)(\alpha I + H)^{-1}\|_2 = \max_{\lambda_i \in \lambda(H)} \left| \frac{\alpha - \lambda_i}{\alpha + \lambda_i} \right| \equiv \sigma(\alpha).$$

由于 $\lambda_i > 0 (i = 1, 2, \dots, n)$ 及 $\alpha > 0$, 容易推得 $\rho(M(\alpha)) \leq \sigma(\alpha) < 1$. 证毕. \square

定理 3.27 表明 HSS 迭代法的收敛速度以 $\sigma(\alpha)$ 为上界, 而 $\sigma(\alpha)$ 仅依赖于 A 的 Hermite 部分 H 的谱 (而不依赖于反 Hermite 部分 S 的谱或系数矩阵 A 的谱).

现引入向量范数 $\|x\|_* = \|(\alpha I + S)x\|_2 (\forall x \in \mathbb{C}^n)$ 且其诱导矩阵范数表示为

$$\|A\|_* = \|(\alpha I + S)A(\alpha I + S)^{-1}\|_2 (\forall A \in \mathbb{C}^{n \times n}),$$

则由定理 3.27 的证明可见

$$\|M(\alpha)\|_* = \|(\alpha I - H)(\alpha I + H)^{-1}(\alpha I - S)(\alpha I + S)^{-1}\|_2 \leq \sigma(\alpha),$$

从而

$$\|x^{(k+1)} - x^*\|_* \leq \sigma(\alpha) \|x^{(k)} - x^*\|_*, \quad k = 0, 1, 2, \dots$$

因此, $\sigma(\alpha)$ 也是 HSS 迭代压缩因子在 $\|\cdot\|_*$ 范数意义下的一个上界. 如果 Hermite 部分 H 的特征值的上界和下界是已知的, 则可得 $\sigma(\alpha)$ (或 $\rho(M(\alpha))$ 和 $\|M(\alpha)\|_*$ 之上界) 的最优参数. 此事实下面的推论中陈述.

推论 3.2 设 $A \in \mathbb{C}^{n \times n}$ 为正定矩阵, $H = \frac{1}{2}(A + A^H)$ 和 $S = \frac{1}{2}(A - A^H)$ 为其 Hermite 和反 Hermite 部分, 且 λ_{\min} 和 λ_{\max} 分别为矩阵 H 的最小和最大特征值, α 是一个正常数. 则

$$\alpha^* \equiv \arg \min_{\alpha} \left\{ \max_{\lambda_{\min} \leq \lambda \leq \lambda_{\max}} \left| \frac{\alpha - \lambda}{\alpha + \lambda} \right| \right\} = \sqrt{\lambda_{\min} \lambda_{\max}},$$

且

$$\sigma(\alpha^*) = \frac{\sqrt{\lambda_{\max}} - \sqrt{\lambda_{\min}}}{\sqrt{\lambda_{\max}} + \sqrt{\lambda_{\min}}} = \frac{\sqrt{\kappa(H)} - 1}{\sqrt{\kappa(H)} + 1},$$

式中: $\kappa(H) = \lambda_{\max}/\lambda_{\min}$ 为 H 的谱条件数.

证明 注意到对任意的 $\alpha > 0$, 函数 $f(\lambda) = \frac{\alpha - \lambda}{\alpha + \lambda}$ 关于 λ 是单调递减的 ($f'(\lambda) = -2\alpha/(\alpha + \lambda)^2 < 0$), 故有

$$\sigma(\alpha) = \max \left\{ \left| \frac{\alpha - \lambda_{\min}}{\alpha + \lambda_{\min}} \right|, \left| \frac{\alpha - \lambda_{\max}}{\alpha + \lambda_{\max}} \right| \right\}. \quad (3.79)$$

若 α^* 是 $\sigma(\alpha)$ 的极小点, 则必有 $\alpha^* = \lambda_{\min} > 0$, $\alpha^* - \lambda_{\max} < 0$,

$$\frac{\alpha^* - \lambda_{\min}}{\alpha^* + \lambda_{\min}} = -\frac{\alpha^* - \lambda_{\max}}{\alpha^* + \lambda_{\max}}.$$

从上式解得

$$\alpha^* = \sqrt{\lambda_{\min} \lambda_{\max}},$$

从而推论的结论成立. 证毕. \square

这里强调在推论 3.2 中, 最优参数 α^* 极小化迭代矩阵谱半径的上界 $\sigma(\alpha)$, 而不是极小化迭代矩阵的谱半径. 推论表明当使用最优参数 α^* 时, HSS 迭代法的收敛速度上界与共轭梯度方法的大致相同, 且当 A 为 Hermite 时, 它们是相同的. 应该指出, 当系数矩阵 A 为正规矩阵时, 有 $HS = SH$, 因此 $\rho(M(\alpha)) = \|M(\alpha)\|_* = \sigma(\alpha)$, 此时最优参数 α^* 极小化所有这三个量.

下面的定理在更一般的情况下分析 IHSS 迭代法. 特别地, 考虑对两步分裂技术的非精确迭代 (比较引理 3.1). 为此, 将 $\|\cdot\|_*$ 推广到 $\|\cdot\|_{M_2}$, 现定义为 $\|x\|_{M_2} = \|M_2 x\|_2 (\forall x \in \mathbb{C}^n)$ 且它诱导矩阵范数为 $\|A\|_{M_2} = \|M_2 A M_2^{-1}\|_2 (\forall A \in \mathbb{C}^{n \times n})$.

定理 3.28 设 $A \in \mathbb{C}^{n \times n}$ 且 $A = M_i - N_i$ ($i = 1, 2$) 为 A 的两个分裂. 如果 $\{x^{(k)}\}$ 是一个如下定义的迭代序列

$$x^{(k+\frac{1}{2})} = x^{(k)} + z^{(k)}, \quad M_1 z^{(k)} = r^{(k)} - p^{(k)} \quad (3.80)$$

满足 $\|p^{(k)}\|_2 \leq \varepsilon_k \|r^{(k)}\|_2$, 其中 $r^{(k)} = b - Ax^{(k)}$, 且

$$x^{(k+1)} = x^{(k+\frac{1}{2})} + z^{(k+\frac{1}{2})}, \quad M_2 z^{(k+\frac{1}{2})} = r^{(k+\frac{1}{2})} - q^{(k+\frac{1}{2})} \quad (3.81)$$

满足 $\|q^{(k+\frac{1}{2})}\|_2 \leq \eta_k \|r^{(k+\frac{1}{2})}\|_2$, 其中 $r^{(k+\frac{1}{2})} = b - Ax^{(k+\frac{1}{2})}$, 则 $\{x^{(k)}\}$ 形如

$$x^{(k+1)} = M_2^{-1} N_2 M_1^{-1} N_1 x^{(k)} + M_2^{-1} (I + N_2 M_1^{-1}) b - M_2^{-1} (N_2 M_1^{-1} p^{(k)} + q^{(k+\frac{1}{2})}). \quad (3.82)$$

进一步, 如果 $x^* \in \mathbb{C}^n$ 是方程组 (3.73) 的精确解, 则有

$$\|x^{(k+1)} - x^*\|_{M_2} \leq (\sigma + \mu \theta \varepsilon_k + \theta(\rho + \theta \nu \varepsilon_k) \eta_k) \|x^{(k)} - x^*\|_{M_2}, \quad k = 0, 1, 2, \dots, \quad (3.83)$$

式中:

$$\begin{aligned} \sigma &= \|N_2 M_1^{-1} N_1 M_2^{-1}\|_2, \quad \rho = \|M_2 M_1^{-1} N_1 M_2^{-1}\|_2, \\ \mu &= \|N_2 M_1^{-1}\|_2, \quad \theta = \|A M_2^{-1}\|_2, \quad \nu = \|M_2 M_1^{-1}\|_2. \end{aligned}$$

特别地, 如果

$$\sigma + \mu \theta \varepsilon_{\max} + \theta(\rho + \theta \nu \varepsilon_{\max}) \eta_{\max} < 1,$$

则迭代序列 $\{x^{(k)}\}$ 收敛到 $x^* \in \mathbb{C}^n$, 其中 $\varepsilon_{\max} = \max_k \{\varepsilon_k\}$ 且 $\eta_{\max} = \max_k \{\eta_k\}$.

证明 由式 (3.80) 并将 $r^{(k)}$ 代入其中, 得

$$x^{(k+\frac{1}{2})} = x^{(k)} + M_1^{-1} (r^{(k)} - p^{(k)}) = M_1^{-1} N_1 x^{(k)} + M_1^{-1} b - M_1^{-1} p^{(k)}. \quad (3.84)$$

类似地, 由式 (3.81) 并将 $r^{(k+\frac{1}{2})}$ 代入其中, 得

$$x^{(k+1)} = x^{(k+\frac{1}{2})} + M_2^{-1} (r^{(k+\frac{1}{2})} - q^{(k+\frac{1}{2})}) = M_2^{-1} N_2 x^{(k+\frac{1}{2})} + M_2^{-1} b - M_2^{-1} q^{(k+\frac{1}{2})}.$$

因此有

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \mathbf{M}_2^{-1} \mathbf{N}_2 (\mathbf{M}_1^{-1} \mathbf{N}_1 \mathbf{x}^{(k)} + \mathbf{M}_1^{-1} \mathbf{b} - \mathbf{M}_1^{-1} \mathbf{p}^{(k)}) + \mathbf{M}_2^{-1} \mathbf{b} - \mathbf{M}_2^{-1} \mathbf{q}^{(k+\frac{1}{2})} \\ &= \mathbf{M}_2^{-1} \mathbf{N}_2 \mathbf{M}_1^{-1} \mathbf{N}_1 \mathbf{x}^{(k)} + \mathbf{M}_2^{-1} (\mathbf{I} + \mathbf{N}_2 \mathbf{M}_1^{-1}) \mathbf{b} - \mathbf{M}_2^{-1} (\mathbf{N}_2 \mathbf{M}_1^{-1} \mathbf{p}^{(k)} + \mathbf{q}^{(k+\frac{1}{2})}). \end{aligned} \quad (3.85)$$

由于 $\mathbf{x}^* \in \mathbb{C}^n$ 是方程组 (3.73) 的精确解, 必满足

$$\mathbf{x}^* = \mathbf{M}_1^{-1} \mathbf{N}_1 \mathbf{x}^* + \mathbf{M}_1^{-1} \mathbf{b}, \quad (3.86)$$

$$\mathbf{x}^* = \mathbf{M}_2^{-1} \mathbf{N}_2 \mathbf{M}_1^{-1} \mathbf{N}_1 \mathbf{x}^* + \mathbf{M}_2^{-1} (\mathbf{I} + \mathbf{N}_2 \mathbf{M}_1^{-1}) \mathbf{b}. \quad (3.87)$$

分别由式 (3.84) 减去式 (3.86), 由式 (3.85) 减去式 (3.87), 有

$$\mathbf{x}^{(k+\frac{1}{2})} - \mathbf{x}^* = \mathbf{M}_1^{-1} \mathbf{N}_1 (\mathbf{x}^{(k)} - \mathbf{x}^*) - \mathbf{M}_1^{-1} \mathbf{p}^{(k)}, \quad (3.88)$$

$$\mathbf{x}^{(k+1)} - \mathbf{x}^* = \mathbf{M}_2^{-1} \mathbf{N}_2 \mathbf{M}_1^{-1} \mathbf{N}_1 (\mathbf{x}^{(k)} - \mathbf{x}^*) - \mathbf{M}_2^{-1} (\mathbf{N}_2 \mathbf{M}_1^{-1} \mathbf{p}^{(k)} + \mathbf{q}^{(k+\frac{1}{2})}). \quad (3.89)$$

式 (3.88) 和式 (3.89) 两边取范数, 得

$$\begin{aligned} \|\mathbf{x}^{(k+\frac{1}{2})} - \mathbf{x}^*\|_{\mathbf{M}_2} &\leq \|\mathbf{M}_1^{-1} \mathbf{N}_1 (\mathbf{x}^{(k)} - \mathbf{x}^*)\|_{\mathbf{M}_2} + \|\mathbf{M}_1^{-1} \mathbf{p}^{(k)}\|_{\mathbf{M}_2} \\ &\leq \|\mathbf{M}_1^{-1} \mathbf{N}_1\|_{\mathbf{M}_2} \cdot \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{\mathbf{M}_2} + \|\mathbf{M}_1^{-1} \mathbf{p}^{(k)}\|_{\mathbf{M}_2} \end{aligned} \quad (3.90)$$

$$\begin{aligned} &\leq \|\mathbf{M}_2 \mathbf{M}_1^{-1} \mathbf{N}_1 \mathbf{M}_2^{-1}\|_2 \cdot \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{\mathbf{M}_2} + \|\mathbf{M}_2 \mathbf{M}_1^{-1}\|_2 \cdot \|\mathbf{p}^{(k)}\|_2, \\ \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|_{\mathbf{M}_2} &\leq \|\mathbf{M}_2^{-1} \mathbf{N}_2 \mathbf{M}_1^{-1} \mathbf{N}_1\|_{\mathbf{M}_2} \cdot \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{\mathbf{M}_2} + \|\mathbf{M}_2^{-1} (\mathbf{N}_2 \mathbf{M}_1^{-1} \mathbf{p}^{(k)} + \mathbf{q}^{(k+\frac{1}{2})})\|_{\mathbf{M}_2} \\ &= \|\mathbf{N}_2 \mathbf{M}_1^{-1} \mathbf{N}_1 \mathbf{M}_2^{-1}\|_2 \cdot \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{\mathbf{M}_2} + \|\mathbf{N}_2 \mathbf{M}_1^{-1} \mathbf{p}^{(k)} + \mathbf{q}^{(k+\frac{1}{2})}\|_2 \\ &\leq \|\mathbf{N}_2 \mathbf{M}_1^{-1} \mathbf{N}_1 \mathbf{M}_2^{-1}\|_2 \cdot \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{\mathbf{M}_2} + \|\mathbf{N}_2 \mathbf{M}_1^{-1}\|_2 \cdot \|\mathbf{p}^{(k)}\|_2 + \|\mathbf{q}^{(k+\frac{1}{2})}\|_2. \end{aligned} \quad (3.91)$$

注意到

$$\|\mathbf{r}^{(k)}\|_2 = \|\mathbf{b} - \mathbf{A} \mathbf{x}^{(k)}\|_2 = \|\mathbf{A}(\mathbf{x}^* - \mathbf{x}^{(k)})\|_2 \leq \|\mathbf{A} \mathbf{M}_2^{-1}\|_2 \cdot \|\mathbf{x}^* - \mathbf{x}^{(k)}\|_{\mathbf{M}_2},$$

$$\|\mathbf{r}^{(k+\frac{1}{2})}\|_2 = \|\mathbf{b} - \mathbf{A} \mathbf{x}^{(k+\frac{1}{2})}\|_2 = \|\mathbf{A}(\mathbf{x}^* - \mathbf{x}^{(k+\frac{1}{2})})\|_2 \leq \|\mathbf{A} \mathbf{M}_2^{-1}\|_2 \cdot \|\mathbf{x}^* - \mathbf{x}^{(k+\frac{1}{2})}\|_{\mathbf{M}_2},$$

由式 (3.88) 和序列 $\{\mathbf{p}^{(k)}\}$ 和 $\{\mathbf{q}^{(k+\frac{1}{2})}\}$ 的定义, 有

$$\|\mathbf{p}^{(k)}\|_2 \leq \varepsilon_k \|\mathbf{r}^{(k)}\|_2 \leq \varepsilon_k \|\mathbf{A} \mathbf{M}_2^{-1}\|_2 \cdot \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{\mathbf{M}_2}, \quad (3.92)$$

$$\begin{aligned} \|\mathbf{q}^{(k+\frac{1}{2})}\|_2 &\leq \eta_k \|\mathbf{r}^{(k+\frac{1}{2})}\|_2 \\ &\leq \eta_k \|\mathbf{A} \mathbf{M}_2^{-1}\|_2 (\|\mathbf{M}_2 \mathbf{M}_1^{-1} \mathbf{N}_1 \mathbf{M}_2^{-1}\|_2 \cdot \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{\mathbf{M}_2} + \|\mathbf{M}_2 \mathbf{M}_1^{-1}\|_2 \cdot \|\mathbf{p}^{(k)}\|_2) \\ &\leq \eta_k \|\mathbf{A} \mathbf{M}_2^{-1}\|_2 (\|\mathbf{M}_2 \mathbf{M}_1^{-1} \mathbf{N}_1 \mathbf{M}_2^{-1}\|_2 + \varepsilon_k \|\mathbf{M}_2 \mathbf{M}_1^{-1}\|_2 \|\mathbf{A} \mathbf{M}_2^{-1}\|_2) \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{\mathbf{M}_2}. \end{aligned} \quad (3.93)$$

将式 (3.92) 和式 (3.93) 代入式 (3.91), 得

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|_{\mathbf{M}_2} \leq [\|\mathbf{N}_2 \mathbf{M}_1^{-1} \mathbf{N}_1 \mathbf{M}_2^{-1}\|_2 + \varepsilon_k \|\mathbf{N}_2 \mathbf{M}_1^{-1}\|_2 \cdot \|\mathbf{A} \mathbf{M}_2^{-1}\|_2 + \eta_k \|\mathbf{A} \mathbf{M}_2^{-1}\|_2]$$

$$\begin{aligned} & \times (\|M_2 M_1^{-1} N_1 M_2^{-1}\| + \varepsilon_k \|M_2 M_1^{-1}\| \cdot \|A M_2^{-1}\|) \cdot \|x^{(k)} - x^*\|_{M_2} \\ & \leq [\sigma + \mu \theta \varepsilon_k + \theta(\rho + \theta \nu \varepsilon_k) \eta_k] \cdot \|x^{(k)} - x^*\|_{M_2}. \end{aligned}$$

证毕. \square

注 3.3 如果在某些应用中可精确求解内方程组, 相应的量 $\{\varepsilon_k\}$ 和 $\{\eta_k\}$ 及 ε_{\max} 和 η_{\max} 等于零, 从而得到 IHSS 迭代法的收敛速度与 HSS 迭代法的相同.

将定理 3.28 应用到 Hermite 和反 Hermite 分裂

$$\begin{aligned} A &= M_1 - N_1 \equiv (\alpha I + H) - (\alpha I - S) \\ &= M_2 - N_2 \equiv (\alpha I + S) - (\alpha I - H), \end{aligned}$$

直接得到下面关于 IHSS 迭代法的收敛性定理.

定理 3.29 设 $A \in \mathbb{C}^{n \times n}$ 为正定矩阵, $H = \frac{1}{2}(A + A^H)$ 和 $S = \frac{1}{2}(A - A^H)$ 为其 Hermite 和反 Hermite 部分, 且 α 是一个正常数. 若 $\{x^{(k)}\}$ 是由 IHSS 迭代法产生的迭代序列且方程组 (3.73) 的精确解是 $x^* \in \mathbb{C}^n$, 则成立

$$\|x^{(k+1)} - x^*\|_* \leq (\sigma(\alpha) + \theta \rho \eta_k)(1 + \theta \varepsilon_k) \|x^{(k)} - x^*\|_*, \quad k = 0, 1, 2, \dots,$$

式中: 范数 $\|\cdot\|_*$ 的定义如前,

$$\rho = \|(\alpha I + S)(\alpha I + H)^{-1}\|_2, \quad \theta = \|A(\alpha I + S)^{-1}\|_2. \quad (3.94)$$

特别地, 若 $(\sigma(\alpha) + \theta \rho \eta_{\max})(1 + \theta \varepsilon_{\max}) < 1$, 则迭代序列 $\{x^{(k)}\}$ 收敛到 $x^* \in \mathbb{C}^n$, 其中 $\varepsilon_{\max} = \max_k \{\varepsilon_k\}$ 且 $\eta_{\max} = \max_k \{\eta_k\}$.

依据定理 3.28, 可选择阈值来极小化两步分裂迭代法的计算工作量. 注意到不需要当 k 增大时阈值 $\{\varepsilon_k\}$ 和 $\{\eta_k\}$ 趋向于零来得到 IHSS 迭代法的收敛性. 下面的定理给出选择阈值 $\{\varepsilon_k\}$ 和 $\{\eta_k\}$ 的一种方式, 使得两步分裂迭代法的原收敛速度 (比较引理 3.1) 可被渐近地恢复.

定理 3.30 设定理 3.28 的条件成立, $\{\tau_1(k)\}$ 和 $\{\tau_2(k)\}$ 是满足 $\tau_1(k) \geq 1$, $\tau_2(k) \geq 1$ 的非降正数序列, 且 $\limsup_{k \rightarrow \infty} \tau_1(k) = \limsup_{k \rightarrow \infty} \tau_2(k) = +\infty$, δ_1 和 δ_2 为属于 $(0, 1)$ 区间的实常数, 满足

$$\varepsilon_k \leq c_1 \delta_1^{\tau_1(k)}, \quad \eta_k \leq c_2 \delta_2^{\tau_2(k)}, \quad k = 0, 1, 2, \dots, \quad (3.95)$$

式中: c_1 和 c_2 为非负常数, 则有

$$\|x^{(k+1)} - x^*\|_{M_2} \leq (\sqrt{\sigma} + \omega \theta \delta^{\tau(k)})^2 \|x^{(k)} - x^*\|_{M_2}, \quad k = 0, 1, 2, \dots,$$

式中:

$$\tau(k) = \min\{\tau_1(k), \tau_2(k)\}, \quad \delta = \max\{\delta_1, \delta_2\}, \quad \omega = \max\left\{\sqrt{c_1 c_2 \nu}, \frac{1}{2\sqrt{\sigma}}(c_1 \mu + c_2 \rho)\right\}.$$

特别地, 有

$$\limsup_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|_{M_2}}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{M_2}} = \sigma,$$

即非精确两步分裂迭代法的收敛速度渐近地与精确两步分裂迭代法的相同.

证明 由式 (3.83) 和式 (3.95), 得

$$\begin{aligned} \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|_{M_2} &\leq (\sigma + \mu\theta\varepsilon_k + \theta(\rho + \theta\nu\varepsilon_k)\eta_k) \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{M_2} \\ &\leq [\sigma + \mu\theta c_1 \delta_1^{\tau_1(k)} + \theta(\rho + \theta\nu c_1 \delta_1^{\tau_1(k)}) c_2 \delta_2^{\tau_2(k)}] \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{M_2} \\ &\leq [\sigma + \mu\theta c_1 \delta^{\tau(k)} + \theta(\rho + \theta\nu c_1 \delta^{\tau(k)}) c_2 \delta^{\tau(k)}] \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{M_2} \\ &= [\sigma + (c_1\mu + c_2\rho)\theta\delta^{\tau(k)} + c_1 c_2 \nu \theta^2 \delta^{2\tau(k)}] \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{M_2} \\ &\leq (\sigma + 2\omega\sqrt{\sigma}\delta^{\tau(k)} + \omega^2\theta^2\delta^{2\tau(k)}) \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{M_2} \\ &= (\sqrt{\sigma} + \omega\theta\delta^{\tau(k)})^2 \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_{M_2}. \end{aligned}$$

由此即得定理的结论. 证毕. \square

由定理 3.29 和定理 3.30 立即可导出下面 IHSS 迭代法的收敛性结果.

定理 3.31 设定理 3.29 的条件成立, $\{\tau_1(k)\}$ 和 $\{\tau_2(k)\}$ 是满足 $\tau_1(k) \geq 1, \tau_2(k) \geq 1$ 的非降正数序列, 且 $\limsup_{k \rightarrow \infty} \tau_1(k) = \limsup_{k \rightarrow \infty} \tau_2(k) = +\infty$, δ_1 和 δ_2 为属于 $(0, 1)$ 区间的实常数, 满足式 (3.95). 则成立

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|_* \leq (\sqrt{\sigma(\alpha)} + \omega\theta\delta^{\tau(k)})^2 \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_*, \quad k = 0, 1, 2, \dots,$$

式中: ρ 和 θ 由式 (3.94) 定义; $\tau(k)$ 和 δ 由式 (3.30) 定义, 且

$$\omega = \max \left\{ \sqrt{c_1 c_2 \rho}, \frac{1}{2\sqrt{\sigma(\alpha)}} (c_1 \sigma(\alpha) + c_2 \rho) \right\}.$$

特别地, 有

$$\limsup_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|_*}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_*} = \sigma(\alpha),$$

即 IHSS 迭代法的收敛速度渐近地与 HSS 迭代法的相同. 此处的范数 $\|\cdot\|_*$ 定义如前.

定理 3.31 表明如果阈值 $\{\varepsilon_k\}$ 和 $\{\eta_k\}$ 如式 (3.95) 中所选择, 则 IHSS 迭代法收敛到方程组 (3.73) 的唯一解 $\mathbf{x}^* \in \mathbb{C}^n$, 且 IHSS 迭代渐近收敛因子的上界趋向于 HSS 迭代的 $\sigma(\alpha)$ (见定理 3.27). 进一步, 式 (3.95) 也可用 $\{\varepsilon_k\}$ 和 $\{\eta_k\}$ 趋向于零来替代.

下面给出 IHSS 迭代法的运算量和存储量估计.

用 ν 表示一次矩阵向量乘积 $\mathbf{A}\mathbf{x}$ 所需要的运算量, $\chi_k(\mathbf{H})$ 和 $\chi_k(\mathbf{S})$ 分别是用阈值 $\{\varepsilon_k\}$ 和 $\{\eta_k\}$ 来非精确求解涉及 \mathbf{H} 和 \mathbf{S} 的内方程组所需的运算量, 则一次 IHSS 迭代的计算工作量可由表 ?? 来估计. 直接计算表明计算 IHSS 迭代每一步的总工作量是 $O(4n - 2\nu + \chi_k(\mathbf{H}) + \chi_k(\mathbf{S}))$.

另外, 简单的计算表明所需要的存储量是储存 $\mathbf{x}^{(k)}, \mathbf{b}, \mathbf{r}^{(k)}, \mathbf{z}^{(k)}$. 对于内方程组的非精确求解器, 只需要另外一些辅助向量, 例如, CG 方法需要大约五个向量 (见第 4 章). 另外, 不需要显式地存储 \mathbf{H} 和 \mathbf{S} , 所需要的仅是执行关于这两个矩阵的矩阵向量乘积子程序. 因此, 所需的存储总量是 $O(n)$, 即与未知量个数的量同阶.

表 3.2 一次 IHSS 迭代所需的工作量

运算	工作量	运算	工作量
$\mathbf{r}^{(k)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}$	$\bar{n} + \nu$	$(\alpha\mathbf{I} + \mathbf{H})\mathbf{z}^{(k+\frac{1}{2})} = \mathbf{r}^{(k)}$	$\chi_k(\mathbf{H})$
$\mathbf{x}^{(k+\frac{1}{2})} = \mathbf{x}^{(k)} + \mathbf{z}^{(k+\frac{1}{2})}$	n	$\mathbf{r}^{(k+\frac{1}{2})} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k+\frac{1}{2})}$	$n + \nu$
$(\alpha\mathbf{I} + \mathbf{S})\mathbf{z}^{(k+\frac{1}{2})} = \mathbf{r}^{(k+\frac{1}{2})}$	$\chi_k(\mathbf{S})$	$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k+\frac{1}{2})} + \mathbf{z}^{(k+1)}$	n

3.4.2 PHSS 迭代法

预处理技术可以大大加快迭代法的收敛速度. 利用 HSS 迭代技术, 文献 [22] 对 2×2 块结构的正半定线性方程组, 建立了一类预处理 HSS 迭代法, 简称 PHSS 迭代法. 理论分析表明方法无条件地收敛到线性方程组的唯一解. 此外, 推导出其收缩因子的上界并给出所涉及迭代参数的最佳选择.

考虑具有块系数矩阵的线性方程组

$$\mathbf{A}\mathbf{x} := \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ -\mathbf{B}^H & \mathbf{O} \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{z} \end{bmatrix} = \begin{bmatrix} \mathbf{f} \\ \mathbf{g} \end{bmatrix} := \mathbf{b}, \quad (3.96)$$

式中: 子矩阵 $\mathbf{A} \in \mathbb{C}^{m \times m}$ 是 Hermite 正定的; $\mathbf{B} \in \mathbb{C}^{m \times n}$ ($m \geq n$) 是列满秩的. 因此, 矩阵 $\mathbf{A} \in \mathbb{C}^{(m+n) \times (m+n)}$ 是非奇异、非 Hermite 半正定矩阵.

块线性方程组 (3.96) 相应于线性约束二次规划问题或鞍点问题的 Kuhn-Tucker 条件. 这些方程组典型地来自于二阶椭圆方程、弹性问题或 Stokes 方程的混合有限元近似. 下面首先给出预处理块线性方程组 (3.96) 一个等价的形式, 然后应用 HSS 方法于预处理块线性方程组, 从而对非 Hermite 半正定方程组 (3.96) 建立一类 PHSS 迭代法.

引入预处理矩阵

$$\mathbf{P} = \begin{bmatrix} \mathbf{A} & \mathbf{O} \\ \mathbf{O} & \mathbf{S} \end{bmatrix} \in \mathbb{C}^{(m+n) \times (m+n)} \text{ 和 } \hat{\mathbf{B}} = \mathbf{A}^{-\frac{1}{2}}\mathbf{B}\mathbf{S}^{-\frac{1}{2}} \in \mathbb{C}^{m \times n}, \quad (3.97)$$

式中: $\mathbf{S} \in \mathbb{C}^{n \times n}$ 为可自由选择的 Hermite 正定子矩阵, 其最佳选择在后面讨论. 定义

$$\begin{aligned} \hat{\mathbf{A}} &= \mathbf{P}^{-\frac{1}{2}}\mathbf{A}\mathbf{P}^{-\frac{1}{2}} = \begin{bmatrix} \mathbf{I} & \hat{\mathbf{B}} \\ -\hat{\mathbf{B}}^H & \mathbf{O} \end{bmatrix}, \quad \begin{bmatrix} \hat{\mathbf{y}} \\ \hat{\mathbf{z}} \end{bmatrix} = \mathbf{P}^{\frac{1}{2}} \begin{bmatrix} \mathbf{y} \\ \mathbf{z} \end{bmatrix} = \begin{bmatrix} \mathbf{A}^{\frac{1}{2}}\mathbf{y} \\ \mathbf{S}^{\frac{1}{2}}\mathbf{z} \end{bmatrix}, \\ \hat{\mathbf{b}} &= \begin{bmatrix} \hat{\mathbf{f}} \\ \hat{\mathbf{g}} \end{bmatrix} = \mathbf{P}^{-\frac{1}{2}}\mathbf{b} = \begin{bmatrix} \mathbf{A}^{-\frac{1}{2}}\mathbf{f} \\ \mathbf{S}^{-\frac{1}{2}}\mathbf{g} \end{bmatrix}. \end{aligned}$$

则线性方程组 (3.96) 可转换成下面等价的预处理形式:

$$\hat{\mathbf{A}} \begin{bmatrix} \hat{\mathbf{y}} \\ \hat{\mathbf{z}} \end{bmatrix} = \hat{\mathbf{b}}. \quad (3.98)$$

显然, 矩阵 $\hat{\mathbf{A}} \in \mathbb{C}^{(m+n) \times (m+n)}$ 的 Hermite 和反 Hermite 部分分别为

$$\hat{\mathbf{H}} = \frac{1}{2}(\hat{\mathbf{A}} + \hat{\mathbf{A}}^H) = \begin{bmatrix} \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} \end{bmatrix}, \quad \hat{\mathbf{S}} = \frac{1}{2}(\hat{\mathbf{A}} - \hat{\mathbf{A}}^H) = \begin{bmatrix} \mathbf{O} & \hat{\mathbf{B}} \\ -\hat{\mathbf{B}}^H & \mathbf{O} \end{bmatrix}.$$

对式 (3.98) 直接应用 HSS 迭代格式, 得

$$\begin{cases} (\alpha I + \widehat{H}) \begin{bmatrix} \widehat{y}^{(k+\frac{1}{2})} \\ \widehat{z}^{(k+\frac{1}{2})} \end{bmatrix} = (\alpha I - \widehat{S}) \begin{bmatrix} \widehat{y}^{(k)} \\ \widehat{z}^{(k)} \end{bmatrix} + \widehat{b}, \\ (\alpha I + \widehat{S}) \begin{bmatrix} \widehat{y}^{(k+1)} \\ \widehat{z}^{(k+1)} \end{bmatrix} = (\alpha I - \widehat{H}) \begin{bmatrix} \widehat{y}^{(k+\frac{1}{2})} \\ \widehat{z}^{(k+\frac{1}{2})} \end{bmatrix} + \widehat{b}. \end{cases}$$

或等价地, 有

$$\begin{bmatrix} \alpha I & \widehat{B} \\ -\widehat{B}^H & \alpha I \end{bmatrix} \begin{bmatrix} \widehat{y}^{(k+1)} \\ \widehat{z}^{(k+1)} \end{bmatrix} = \begin{bmatrix} \frac{\alpha(\alpha-1)}{\alpha+1} I & -\frac{\alpha-1}{\alpha+1} \widehat{B} \\ \widehat{B}^H & \alpha I \end{bmatrix} \begin{bmatrix} \widehat{y}^{(k)} \\ \widehat{z}^{(k)} \end{bmatrix} + \begin{bmatrix} \frac{2\alpha}{\alpha+1} \widehat{f} \\ 2\widehat{g} \end{bmatrix}.$$

从而按原变量, 得

$$\begin{bmatrix} \alpha A & B \\ -B^H & \alpha S \end{bmatrix} \begin{bmatrix} y^{(k+1)} \\ z^{(k+1)} \end{bmatrix} = \begin{bmatrix} \frac{\alpha(\alpha-1)}{\alpha+1} A & -\frac{\alpha-1}{\alpha+1} B \\ B^H & \alpha S \end{bmatrix} \begin{bmatrix} y^{(k)} \\ z^{(k)} \end{bmatrix} + \begin{bmatrix} \frac{2\alpha}{\alpha+1} f \\ 2g \end{bmatrix}, \quad (3.99)$$

这导致下面求解块线性方程组 (3.96) 的 PHSS 迭代法.

算法 3.8 (PHSS 迭代法)

步 1, 给定一个初始向量 $x^{(0)} = \begin{bmatrix} y^{(0)} \\ z^{(0)} \end{bmatrix} \in \mathbb{C}^{m+n}$, 参数 $\alpha > 0$, 容许误差 $\varepsilon > 0$.

置 $k := 0$.

步 2, 用 HSS 迭代法求解式 (3.99) 计算 $x^{(k+1)} = \begin{bmatrix} y^{(k+1)} \\ z^{(k+1)} \end{bmatrix} \in \mathbb{C}^{m+n}$.

步 3, 计算残差 $r^{(k+1)} = b - Ax^{(k+1)}$. 若 $\|r^{(k+1)}\|_2 \leq \varepsilon$, 停算. 否则, 置 $k := k+1$, 转步 2.

明显地, PHSS 迭代法可等价地重写为

$$\begin{bmatrix} y^{(k+1)} \\ z^{(k+1)} \end{bmatrix} = \mathcal{T}(\alpha) \begin{bmatrix} y^{(k)} \\ z^{(k)} \end{bmatrix} + \mathcal{M}(\alpha)^{-1} \begin{bmatrix} f \\ g \end{bmatrix}, \quad (3.100)$$

式中:

$$\begin{cases} \mathcal{T}(\alpha) = \begin{bmatrix} \alpha A & B \\ -B^H & \alpha S \end{bmatrix}^{-1} \begin{bmatrix} \frac{\alpha(\alpha-1)}{\alpha+1} A & -\frac{\alpha-1}{\alpha+1} B \\ B^H & \alpha S \end{bmatrix}, \\ \mathcal{M}(\alpha)^{-1} = \begin{bmatrix} \alpha A & B \\ -B^H & \alpha S \end{bmatrix}^{-1} \begin{bmatrix} \frac{2\alpha}{\alpha+1} I & O \\ O & 2I \end{bmatrix}. \end{cases} \quad (3.101)$$

这里 $\mathcal{T}(\alpha)$ 是 PHSS 迭代法的迭代矩阵. 事实上, 式 (3.100) 也可来自于系数矩阵 A 的分裂

$$A = \mathcal{M}(\alpha) - \mathcal{N}(\alpha) = \begin{bmatrix} \frac{\alpha+1}{2} A & \frac{\alpha+1}{2\alpha} B \\ -\frac{1}{2} B^H & \frac{\alpha}{2} S \end{bmatrix} - \begin{bmatrix} \frac{\alpha-1}{2} A & -\frac{\alpha-1}{2\alpha} B \\ \frac{1}{2} B^H & \frac{\alpha}{2} S \end{bmatrix}. \quad (3.102)$$

在实际计算中, PHSS 迭代法的每次迭代需要求解具有如下系数矩阵的线性方程组

$$\widetilde{\mathcal{M}}(\alpha) = \begin{bmatrix} \alpha A & B \\ -B^H & \alpha S \end{bmatrix}, \quad (3.103)$$

或等价地, 有

$$\mathcal{M}(\alpha) = \begin{bmatrix} \frac{\alpha+1}{2}A & \frac{\alpha+1}{2\alpha}B \\ -\frac{1}{2}B^H & \frac{\alpha}{2}S \end{bmatrix}. \quad (3.104)$$

由于这些矩阵是正定的, 可用另一个迭代过程非精确地求解前述线性方程组, 如 HSS 迭代. 这导致非 Hermite 正半定线性方程组 (3.96) 的一个非精确预处理 Hermite-反 Hermite 分裂迭代法 (简记为 IPHSS), 并已经在前面有所讨论.

为了对 PHSS 迭代法进行收敛性分析, 通过简单推导、奇异值分解等可得式 (3.101) 中矩阵 $\mathcal{T}(\alpha)$ 的显式表达式、特征值结构以及谱半径, 这由下面的三个引理描述, 其证明详见文献 [22].

引理 3.2 考虑线性方程组 (3.96). 设 $A \in \mathbb{C}^{m \times m}$ 是 Hermite 正定的, $B \in \mathbb{C}^{m \times n}$ 列满秩, $\alpha > 0$ 是一个给定常数. 假设 $S \in \mathbb{C}^{n \times n}$ 是 Hermite 正定矩阵. 则划分式 (3.101) 中的迭代矩阵 $\mathcal{T}(\alpha)$ 为

$$\mathcal{T}(\alpha) = \begin{bmatrix} \mathcal{T}_{11}(\alpha) & \mathcal{T}_{12}(\alpha) \\ \mathcal{T}_{21}(\alpha) & \mathcal{T}_{22}(\alpha) \end{bmatrix},$$

式中:

$$\begin{aligned} \mathcal{T}_{11}(\alpha) &= \frac{\alpha-1}{\alpha+1}I - \frac{2}{\alpha+1}A^{-1}B\widetilde{S}(\alpha)^{-1}B^H, & \mathcal{T}_{12}(\alpha) &= -\frac{2\alpha}{\alpha+1}A^{-1}B\widetilde{S}(\alpha)^{-1}S, \\ \mathcal{T}_{21}(\alpha) &= \frac{2\alpha}{\alpha+1}S\widetilde{S}(\alpha)^{-1}B^HA^{-1}, & \mathcal{T}_{22}(\alpha) &= -\frac{\alpha-1}{\alpha+1}I + \frac{2\alpha^2}{\alpha+1}\widetilde{S}(\alpha)^{-1}S, \end{aligned}$$

且

$$\widetilde{S}(\alpha) = \alpha S + \frac{1}{\alpha}B^HA^{-1}B$$

为式 (3.103) 中矩阵 $\widetilde{\mathcal{M}}(\alpha)$ 的 Schur 补.

引理 3.3 设引理 3.2 中的条件满足. 如果 $\sigma_k (k=1, 2, \dots, n)$ 是式 (3.97) 中矩阵 $\widehat{B} \in \mathbb{C}^{m \times n}$ 的正奇异值, 则 PHSS 方法迭代矩阵 $\mathcal{T}(\alpha)$ 的特征值是具有重数 $m-n$ 的

$$\frac{\alpha-1}{\alpha+1}$$

和

$$\frac{1}{(\alpha+1)(\alpha^2+\sigma_k^2)} \left(\alpha(\alpha^2-\sigma_k^2) \pm \sqrt{(\alpha^2+\sigma_k^2)^2-4\alpha^4\sigma_k^2} \right), \quad k=1, 2, \dots, n.$$

引理 3.4 设引理 3.2 中的条件满足. 如果 $\sigma_k (k=1, 2, \dots, n)$ 是式 (3.97) 中矩阵 $\hat{B} \in \mathbb{C}^{m \times n}$ 的正奇异值, 且 λ 是 PHSS 方法迭代矩阵 $\mathcal{T}(\alpha)$ 的主特征值, 即 $\mathcal{T}(\alpha)$ 的谱半径可由 $|\lambda|$ 达到, 则对 $k=1, 2, \dots, n$, 成立

$$|\lambda| = \begin{cases} \frac{|\alpha-1|}{\alpha+1}, & \\ \frac{\alpha}{\alpha+1} \left(\frac{|\alpha^2 - \sigma_k^2|}{\alpha^2 + \sigma_k^2} + \sqrt{\frac{1}{\alpha^2} - \frac{4\alpha^2 \sigma_k^2}{(\alpha^2 + \sigma_k^2)^2}} \right), & \text{对 } \alpha^2 + \sigma_k^2 > 2\alpha^2 \sigma_k; \\ \sqrt{\frac{\alpha-1}{\alpha+1}}, & \text{对 } \alpha^2 + \sigma_k^2 \leq 2\alpha^2 \sigma_k. \end{cases}$$

基于引理 3.4, 下面证明求解线性方程组 (3.96) 的 PHSS 迭代法的收敛性.

定理 3.32 考虑线性方程组 (3.96). 设 $A \in \mathbb{C}^{m \times m}$ 是 Hermite 正定的, $B \in \mathbb{C}^{m \times n}$ 列满秩, $\alpha > 0$ 是一个给定常数. 假设 $S \in \mathbb{C}^{n \times n}$ 是 Hermite 正定矩阵, 则

$$\rho(\mathcal{T}(\alpha)) < 1, \quad \forall \alpha > 0,$$

即 PHSS 迭代序列收敛到线性方程组 (3.96) 的精确解.

证明 显然有

$$\frac{|\alpha-1|}{\alpha+1} < 1 \quad (\forall \alpha > 0) \quad \text{和} \quad \sqrt{\frac{\alpha-1}{\alpha+1}} < 1 \quad (\forall \alpha > 1).$$

因为对 $k=1, 2, \dots, n$, 当 $\alpha^2 + \sigma_k^2 > 2\alpha^2 \sigma_k$ 时, 成立

$$\begin{aligned} & \frac{\alpha}{\alpha+1} \left(\frac{|\alpha^2 - \sigma_k^2|}{\alpha^2 + \sigma_k^2} + \sqrt{\frac{1}{\alpha^2} - \frac{4\alpha^2 \sigma_k^2}{(\alpha^2 + \sigma_k^2)^2}} \right) \\ & < \frac{\alpha}{\alpha+1} \left(\frac{|\alpha^2 - \sigma_k^2|}{\alpha^2 + \sigma_k^2} + \frac{1}{\alpha} \right) < \frac{\alpha}{\alpha+1} \left(1 + \frac{1}{\alpha} \right) = 1, \end{aligned}$$

利用引理 3.4 易见对 $\forall \alpha > 0$ 均成立 $\rho(\mathcal{T}(\alpha)) < 1$. 证毕. \square

下面的定理描述了 PHSS 迭代法的最佳迭代参数和相应的渐近收敛因子.

定理 3.33 考虑线性方程组 (3.96). 设 $A \in \mathbb{C}^{m \times m}$ 是 Hermite 正定的, $B \in \mathbb{C}^{m \times n}$ 列满秩, $\alpha > 0$ 是一个给定常数. 假设 $S \in \mathbb{C}^{n \times n}$ 是 Hermite 正定矩阵. 如果 $\sigma_k (k=1, 2, \dots, n)$ 是矩阵 $A^{-\frac{1}{2}} B S^{-\frac{1}{2}} \in \mathbb{C}^{m \times n}$ 的正奇异值, 且 $\sigma_{\min} = \min_{1 \leq k \leq n} \sigma_k$ 和 $\sigma_{\max} = \max_{1 \leq k \leq n} \sigma_k$, 则对线性方程组 (3.96) 的 PHSS 迭代法, 迭代参数 α 的最佳值由下式给出,

$$\alpha^* = \arg \min_{\alpha} \rho(\mathcal{T}(\alpha)) = \sqrt{\sigma_{\min} \sigma_{\max}},$$

且相应地, 有

$$\rho(\mathcal{T}(\alpha^*)) = \frac{\sigma_{\max} - \sigma_{\min}}{\sigma_{\max} + \sigma_{\min}}.$$

证明 首先注意到下面两个事实成立:

(1) 当 $\alpha \leq 1$ 时, $\alpha^2 + \sigma_k^2 > 2\alpha^2\sigma_k$, $k = 1, 2, \dots, n$.

(2) 当 $\alpha > 1$ 时, 记 $\alpha_- = \alpha^2 - \alpha\sqrt{\alpha^2 - 1}$ 和 $\alpha_+ = \alpha^2 + \alpha\sqrt{\alpha^2 - 1}$, 则

① $\alpha^2 + \sigma_k^2 > 2\alpha^2\sigma_k$ 当且仅当 $\sigma_k \in (0, \alpha_-) \cup (\alpha_+, +\infty)$, $k \in \{1, 2, \dots, n\}$;

② $\alpha^2 + \sigma_k^2 \leq 2\alpha^2\sigma_k$ 当且仅当 $\sigma_k \in [\alpha_-, \alpha_+]$, $k \in \{1, 2, \dots, n\}$;

③ $\frac{\alpha - 1}{\alpha + 1} < \sqrt{\frac{\alpha - 1}{\alpha + 1}}$.

设

$$\theta(\alpha, \sigma) = \frac{\alpha}{\alpha + 1} \left(\frac{|\alpha^2 - \sigma^2|}{\alpha^2 + \sigma^2} + \sqrt{\frac{1}{\alpha^2} - \frac{4\alpha^2\sigma^2}{(\alpha^2 + \sigma^2)^2}} \right).$$

则基于事实 (1) 和 (2), 根据引理 3.4, 易知

$$\rho(\mathcal{T}(\alpha)) = \begin{cases} \max \left\{ \frac{1 - \alpha}{1 + \alpha}, \max_{1 \leq k \leq n} \theta(\alpha, \sigma_k) \right\}, & \alpha \leq 1, \\ \max \left\{ \sqrt{\frac{\alpha - 1}{\alpha + 1}}, \max_{\substack{\sigma_k < \alpha_- \text{ 或 } \sigma_k > \alpha_+ \\ k \in \{1, 2, \dots, n\}}} \theta(\alpha, \sigma_k) \right\}, & \alpha > 1. \end{cases} \quad (3.105)$$

对任意固定的 $\beta > 0$, 定义两个函数 $\theta_1, \theta_2: (0, +\infty) \rightarrow (0, +\infty)$, 即

$$\theta_1(t) = \frac{\beta - t}{\beta + t}, \quad \theta_2(t) = \frac{1}{\beta} - \frac{4\beta t}{(\beta + t)^2}.$$

直接计算, 得

$$\frac{d\theta_1(t)}{dt} = \frac{-2\beta}{(\beta + t)^2}, \quad \frac{d\theta_2(t)}{dt} = \frac{4\beta(t - \beta)}{(\beta + t)^3}.$$

则

(1) 对 $\alpha \leq 1$ 且 $\sigma_{\min} \leq \alpha \leq \sigma_{\max}$, 有

$$\max_{1 \leq k \leq n} \theta(\alpha, \sigma_k) = \max\{\theta(\alpha, \sigma_{\min}), \theta(\alpha, \sigma_{\max})\}.$$

(2) 对 $\alpha > 1$ 且 $\sigma_{\min} < \alpha_-$ 或 $\sigma_{\max} > \alpha_+$, 有

$$\max_{\substack{\sigma_k < \alpha_- \text{ 或 } \sigma_k > \alpha_+ \\ k \in \{1, 2, \dots, n\}}} \theta(\alpha, \sigma_k) = \max\{\theta(\alpha, \sigma_{\min}), \theta(\alpha, \sigma_{\max})\}.$$

因此, 当 $\alpha \leq 1$ 时, 最佳参数 α^* 必须满足 $\sigma_{\min} \leq \alpha^* \leq \sigma_{\max}$ 和下面三个条件之一:

$$\frac{1 - \alpha^*}{1 + \alpha^*} = \theta(\alpha^*, \sigma_{\min}) \geq \theta(\alpha^*, \sigma_{\max});$$

$$\frac{1 - \alpha^*}{1 + \alpha^*} = \theta(\alpha^*, \sigma_{\max}) \geq \theta(\alpha^*, \sigma_{\min});$$

$$\theta(\alpha^*, \sigma_{\min}) = \theta(\alpha^*, \sigma_{\max}) \geq \frac{1 - \alpha^*}{1 + \alpha^*}.$$

且当 $\alpha > 1$ 时, 最佳参数 α^* 必须满足 $\sigma_{\min} < \alpha_-^*$ 或 $\sigma_{\max} > \alpha_+^*$ 及下面三个条件之一:

$$\sqrt{\frac{\alpha^* - 1}{\alpha^* + 1}} = \theta(\alpha^*, \sigma_{\min}) \geq \theta(\alpha^*, \sigma_{\max});$$

$$\sqrt{\frac{\alpha^* - 1}{\alpha^* + 1}} = \theta(\alpha^*, \sigma_{\max}) \geq \theta(\alpha^*, \sigma_{\min});$$

$$\theta(\alpha^*, \sigma_{\min}) = \theta(\alpha^*, \sigma_{\max}) \geq \sqrt{\frac{\alpha^* - 1}{\alpha^* + 1}}.$$

其中 $\alpha_-^* = (\alpha^*)^2 - \alpha^* \sqrt{(\alpha^*)^2 - 1}$ 和 $\alpha_+^* = (\alpha^*)^2 + \alpha^* \sqrt{(\alpha^*)^2 - 1}$. 直接求解当 $a \leq 1$ 的三个条件和 $a > 1$ 的三个条件, 可得 α^* 的具体表达式. 然后将 α^* 代入式 (3.105), 可得 $\rho(\mathcal{T}(\alpha^*))$. 证毕. \square

由定理 3.33 可见应该选择 Hermite 正定矩阵 $S \in \mathbb{C}^{n \times n}$, 使得具有系数矩阵 S 的线性方程组是易于求解的且矩阵 $A^{-\frac{1}{2}}BS^{-\frac{1}{2}} \in \mathbb{C}^{m \times n}$ 的奇异值是紧密聚集的, 换言之, S 应该是 Schur 补 $B^H A^{-1}B \in \mathbb{C}^{n \times n}$ 的一个好的预处理子. 寻求好的矩阵 S 对加快 PHSS 迭代法的收敛速度是至关重要的. 至于最佳参数 $\alpha^* = \sqrt{\sigma_{\min}\sigma_{\max}}$ 在实际应用中的计算, 由于矩阵 $A^{-\frac{1}{2}}BS^{-\frac{1}{2}}$ 的最小和最大奇异值 σ_{\min} 和 σ_{\max} 分别等于矩阵 $S^{-1}B^H A^{-1}B$ 的最小和最大特征值 λ_{\min} 和 λ_{\max} 的平方根, 即 $\sigma_{\min} = \sqrt{\lambda_{\min}}$ 和 $\sigma_{\max} = \sqrt{\lambda_{\max}}$, 从而通过计算 λ_{\min} 和 λ_{\max} , 可得 $\alpha^* = \sqrt[4]{\lambda_{\min}\lambda_{\max}}$.

3.5 迭代法的加速方法

使用一阶线性定常迭代法 (3.4), 或者其他的迭代格式求解方程组 (3.1) 时, 迭代过程可能收敛, 也可能不收敛. 即便收敛, 也可能收敛得很慢. 无论是哪种情况, 都希望找到一种改进方法, 使得不收敛的格式变得收敛, 收敛慢的格式变得收敛快.

给定第 k 次迭代后的估计解 $x^{(k)}$, 由迭代格式 (3.4) 得到了一个新的估计解 $Bx^{(k)} + f$, 再将它与 $x^{(k)}$, $x^{(k-1)}$ 等进行适当的组合, 就可能得到更快的迭代收敛速度. 这就是迭代法加速的基本思想.

3.5.1 外推方法

在迭代格式 (3.4) 中引入一个参数 $\gamma \neq 0$, 构造新的迭代格式

$$x^{(k+1)} = (1 - \gamma)x^{(k)} + \gamma(Bx^{(k)} + f) := B_\gamma x^{(k)} + \gamma f, \quad k = 0, 1, \dots, \quad (3.106)$$

式中:

$$B_\gamma = (1 - \gamma)I + \gamma B. \quad (3.107)$$

当 $\gamma = 1$ 时, 就退化为原迭代格式 (3.4). 若迭代格式 (3.106) 收敛到某个 x^* , 则有

$$x^* = (1 - \gamma)x^* + \gamma(Bx^* + f) \implies x^* = Bx^* + f.$$

所以, 不管 $\gamma \neq 0$ 如何选取, 当迭代格式 (3.106) 收敛时, 它必定收敛到原来方程组 $\mathbf{x} = \mathbf{B}\mathbf{x} + \mathbf{f}$ 的解. 选取适当的 \mathbf{B} , 根据式 (3.3) 中的定义可知, 这样的解也是 $\mathbf{A}\mathbf{x} = \mathbf{b}$ 的解. 因此, 可望选取一个比较好的参数 γ , 使得迭代格式 (3.106) 收敛得尽可能快. 这就是外推法, 它依赖于原迭代矩阵 \mathbf{B} 的选择. 所以分别有 Jacobi 迭代、Gauss-Seidel 迭代和 SOR 迭代的外推格式.

注 3.4 设 \mathbf{B} 是 Gauss-Seidel 迭代矩阵, 即使 γ 取为 SOR 迭代法的松弛因子 ω , 式 (3.106) 与 SOR 迭代法也是截然不同的. 在式 (3.106) 中, 当 $\mathbf{B}\mathbf{x}^{(k)} + \mathbf{f}$ 的每个分量计算完之后, 才能再与 $\mathbf{x}^{(k)}$ 作加权平均得到 $\mathbf{x}^{(k+1)}$.

为了使得外推格式 (3.106) 收敛得尽可能快, 需要选择参数 γ 使得 $\rho(\mathbf{B}_\gamma)$ 尽可能小. 如果 \mathbf{B} 的所有特征值均为实数, 且分布在 $[a, b]$ 之间, 则有

$$\begin{aligned}\rho(\mathbf{B}_\gamma) &= \max_{\lambda \in \lambda(\mathbf{B})} |(1-\gamma) + \gamma\lambda| \leq \max_{\lambda \in [a, b]} |(1-\gamma) + \gamma\lambda| \\ &= \max\{|(1-\gamma) + \gamma a|, |(1-\gamma) + \gamma b|\} \\ &= \max\{|1 + (a-1)\gamma|, |1 + (b-1)\gamma|\}.\end{aligned}$$

上式第 2 个等式成立是因为 $(1-\gamma) + \gamma\lambda$ 是 λ 的线性函数, 其在 $[a, b]$ 上的最大绝对值必在端点取得. 为使 $\rho(\mathbf{B}_\gamma)$ 尽可能小, 只要右端项尽可能小即可. 下面的定理给出了最佳 γ 的取法.

定理 3.34 设 \mathbf{B} 的所有特征值均为实数, 且分布在 $[a, b]$ 之间, 又假定 $1 \notin [a, b]$, 则函数

$$f(\gamma) = \max\{|1 + (a-1)\gamma|, |1 + (b-1)\gamma|\}$$

在极小点

$$\gamma_{\min} = \frac{2}{2-a-b} \quad (3.108)$$

处取得极小值

$$f_{\min} = \begin{cases} 1 - (a-1)|\gamma_{\min}|, & \text{当 } 1 < a \leq b, \\ 1 - (1-b)|\gamma_{\min}|, & \text{当 } a \leq b < 1. \end{cases}$$

根据定理 3.34, 得

$$\rho(\mathbf{B}_\gamma) = f_{\min} = 1 - d|\gamma_{\min}| < 1, \quad (3.109)$$

式中: d 为 1 到 $[a, b]$ 的距离, 即

$$d = \begin{cases} a-1, & 1 < a \leq b, \\ 1-b, & a \leq b < 1. \end{cases}$$

所以, 取 $\gamma = \gamma_{\min}$ 时, 式 (3.106) 收敛. 由于 $\gamma_{\min} \neq 1$, 故取 $\gamma = \gamma_{\min}$ 时的式 (3.106) 要比式 (3.4) (即取 $\gamma = 1$ 时的迭代法) 收敛得更快. 此外, 为了使得 f_{\min} 尽可能小, 可取 a, b 分别为迭代矩阵 \mathbf{B} 的最小和最大特征值.

3.5.2 整体校正方法

设 x_1, x_2, \dots, x_m ($m > 1$) 是方程组 (3.1) 按照某种方式获取的互异近似解, 并设 $Ax_i \neq b$ ($i = 1, 2, \dots, m$). 利用诸 x_i 提供的信息, 构造向量 x , 使得它比 x_i 都更接近于方程组 (3.1) 的精确解 x^* , 即要求

$$\|b - Ax\|_2 < \min_{1 \leq i \leq m} \|b - Ax_i\|_2, \quad (3.110)$$

称满足式 (3.110) 的向量 x 为方程组 (3.1) 关于 x_1, x_2, \dots, x_m 的校正解, 确定 x 的过程称为校正过程. 建立整体校正模型

$$\begin{cases} x = \sum_{i=1}^m \alpha_i x_i, & \sum_{i=1}^m \alpha_i = 1, \\ \|b - Ax\|_2 = \min. \end{cases} \quad (3.111)$$

所谓整体校正, 是指 x_i 的每一个分量对 x 的对应分量的贡献比例相同. 下面求解模型 (3.111).

任意选定 $s \in \{1, 2, \dots, m\}$, 则有

$$\begin{aligned} x &= x_s + \sum_{i \neq s} \alpha_i (x_i - x_s), \\ r(x) &= b - Ax = (b - Ax_s) - \sum_{i \neq s} \alpha_i (Ax_i - Ax_s) \\ &= r_s + \sum_{i \neq s} \alpha_i (r_i - r_s). \end{aligned} \quad (3.112)$$

记

$$\begin{aligned} \beta_i &= r_i - r_s, \quad Q_s = [\beta_1, \dots, \beta_{s-1}, \beta_{s+1}, \dots, \beta_m], \\ y_s &= (\alpha_1, \dots, \alpha_{s-1}, \alpha_{s+1}, \dots, \alpha_m)^T. \end{aligned}$$

那么式 (3.112) 可写为 $r(x) = r_s + Q_s y_s$, 从而求

$$x = \sum_{i=1}^m \alpha_i x_i = \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_m x_m$$

使得 $\|r(x)\|_2 = \min$, 等价于求 $y_s \in \mathbb{C}^{m-1}$ 使得 $\|r_s + Q_s y_s\|_2 = \min$, 其极小范数解为 (见定理 5.2 和定理 5.3)

$$y_s = -Q_s^\dagger r_s. \quad (3.113)$$

再由 $\alpha_s = 1 - \sum_{i \neq s} \alpha_i$, 即可确定数组 $\alpha_1, \dots, \alpha_m$, 从而可确定校正解 x .

特别地, 对于 $m = 2$, 取 $s = 2$ 时, 有

$$\begin{cases} \alpha_1 = y_2 = -Q_2^\dagger r_2 = -\beta_1^\dagger r_2 = -\frac{(r_1 - r_2)^H r_2}{\|r_1 - r_2\|_2^2}, \\ \alpha_2 = 1 - \alpha_1 = \frac{(r_1 - r_2)^H r_1}{\|r_1 - r_2\|_2^2}. \end{cases} \quad (3.114)$$

取 $s = 1$ 时,

$$\begin{cases} \alpha_2 = y_1 = -Q_1^\dagger r_1 = -\beta_2^\dagger r_1 = -\frac{(r_2 - r_1)^H r_1}{\|r_2 - r_1\|_2^2}, \\ \alpha_1 = 1 - \alpha_2 = \frac{(r_2 - r_1)^H r_2}{\|r_2 - r_1\|_2^2}. \end{cases} \quad (3.115)$$

易见, 式 (3.114) 与式 (3.115) 的结果一致.

下面来看模型 (3.111) 的几何意义 (图 3.1). 把 $r(x)$ 和 r_i 分别看作以坐标原点为

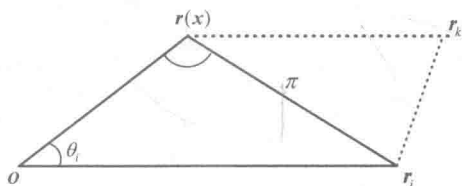


图 3.1 整体校正模型

起点, 以 $r(x)$ 和 r_i 为终点的空间向量, 那么由 $r(x) = \alpha_1 r_1 + \alpha_2 r_2 + \cdots + \alpha_m r_m$ 及 $\alpha_1 + \alpha_2 + \cdots + \alpha_m = 1$ 可得: 点 $r(x)$ 在由点 r_1, r_2, \cdots, r_m 确定的超平面 π 上, 且当 $i \neq s$ 时, 有

$$\begin{aligned} (r_i - r_s)^H r(x) &= \beta_i^H (r_s + Q_s y_s) \\ &= \beta_i^H (r_s - Q_s Q_s^\dagger r_s) = \beta_i^H [r_s - (Q_s Q_s^\dagger)^H r_s] \\ &= [\beta_i^H - (Q_s Q_s^\dagger \beta_i)^H] r_s = (\beta_i^H - \beta_i^H) r_s = 0. \end{aligned} \quad (3.116)$$

上述推导过程用到了广义逆矩阵的性质

$$(Q_s Q_s^\dagger)^H = Q_s Q_s^\dagger, \quad Q_s Q_s^\dagger Q_s = Q_s.$$

式 (3.116) 表明 $r(x) \perp \pi$. 由于“垂足”为 $r(x)$, 所以

$$\|r(x)\|_2 \leq \|r_i\|_2, \quad i = 1, 2, \cdots, m. \quad (3.117)$$

定理 3.35 模型 (3.111) 的校正解 x 满足 $\|r(x)\|_2 < \min_{1 \leq i \leq m} \|r_i\|_2$ 的充分必要条件是 $x \neq x_i (i = 1, 2, \cdots, m)$.

证明 必要性. 由 $\|r(x)\|_2 < \|r_i\|_2$, 得 $r(x) \neq r_i$, 即 $b - Ax \neq b - Ax_i$, 从而有

$$x \neq x_i, \quad (i = 1, 2, \cdots, m).$$

充分性. 由 $x \neq x_i$, 得 $Ax \neq Ax_i$, 即 $r(x) \neq r_i$, 从而有

$$\|r(x)\|_2 \neq \|r_i\|_2 \quad (\text{垂足唯一性}).$$

再由式 (3.117), 得 $\|\mathbf{r}(\mathbf{x})\|_2 < \|\mathbf{r}_i\|_2, i = 1, 2, \dots, m$. 证毕. \square

如果 $m = n + 1$, 且 $\beta_1, \dots, \beta_{s-1}, \beta_{s+1}, \dots, \beta_m$ 线性无关, 则由 $\mathbf{r}(\mathbf{x}) \perp \beta_i (i \neq s)$ 可得 $\mathbf{r}(\mathbf{x}) = \mathbf{0}$, 从而校正解 $\mathbf{x} = \mathbf{x}^*$.

为了衡量校正过程的效果, 引进整体缩减系数

$$\alpha = \frac{\|\mathbf{r}(\mathbf{x})\|_2}{\min_{1 \leq i \leq m} \|\mathbf{r}_i\|_2}. \quad (3.118)$$

由式 (3.117) 知 $0 \leq \alpha \leq 1$. 当 $\alpha = 1$ 时, 称校正过程失败; 当 $0 \leq \alpha < 1$ 时, 称校正过程成功. 特别地, 当 $\alpha = 0$ 时, 称校正过程完成, 此时校正解 $\mathbf{x} = \mathbf{x}^*$. 若用 θ_i 表示 $\mathbf{r}(\mathbf{x})$ 与 \mathbf{r}_i 之间的夹角 (图 3.1), 则由 $\mathbf{r}(\mathbf{x}) \perp \pi$ 及 $\mathbf{r}(\mathbf{x}) \in \pi$ 可得 $\mathbf{r}(\mathbf{x}) \perp (\mathbf{r}(\mathbf{x}) - \mathbf{r}_i)$. 于是有

$$\begin{aligned} \sin \theta_i &= \frac{\|\mathbf{r}(\mathbf{x}) - \mathbf{r}_i\|_2}{\|\mathbf{r}_i\|_2}, \\ \cos \theta_i &= \frac{\|\mathbf{r}(\mathbf{x})\|_2}{\|\mathbf{r}_i\|_2} = \left(1 - \frac{\|\mathbf{r}(\mathbf{x}) - \mathbf{r}_i\|_2^2}{\|\mathbf{r}_i\|_2^2}\right)^{\frac{1}{2}} \\ &= \left(1 - \frac{\|\mathbf{Q}_i \mathbf{y}_i\|_2^2}{\|\mathbf{r}_i\|_2^2}\right)^{\frac{1}{2}} = \left(1 - \frac{\|\mathbf{Q}_i \mathbf{Q}_i^\dagger \mathbf{r}_i\|_2^2}{\|\mathbf{r}_i\|_2^2}\right)^{\frac{1}{2}}. \end{aligned}$$

选取 s_0 满足 $\|\mathbf{r}_{s_0}\|_2 < \min_{1 \leq i \leq m} \|\mathbf{r}_i\|_2$, 则有

$$\alpha = \frac{\|\mathbf{r}(\mathbf{x})\|_2}{\|\mathbf{r}_{s_0}\|_2} = \left(1 - \frac{\|\mathbf{Q}_{s_0} \mathbf{Q}_{s_0}^\dagger \mathbf{r}_{s_0}\|_2^2}{\|\mathbf{r}_{s_0}\|_2^2}\right)^{\frac{1}{2}}. \quad (3.119)$$

应用式 (3.119) 能够事先估计校正过程的效果.

将整体校正过程施加于任何一种求解方程组 (3.1) 的迭代格式, 都可以改善原始式的收敛性态. 设求解方程组 (3.1) 的单步迭代格式为

$$\mathbf{x}^{(i)} = \varphi(\mathbf{x}^{(i-1)}) \quad (i = 1, 2, 3, \dots). \quad (3.120)$$

改造格式 (3.120) 为以下算法.

算法 3.9 (整体校正法)

步 1, 输入矩阵 \mathbf{A} , 右端向量 \mathbf{b} , 初始点 $\mathbf{x}^{(0)}$, 精度要求 ε , 最大迭代次数 N , 取定整数 m 及 $s \in [1, m]$. 置 $k := 1$.

步 2, 置 $\mathbf{x}_1^{(k)} = \mathbf{x}^{(k-1)}$, 计算

$$\mathbf{x}_i^{(k)} = \varphi(\mathbf{x}_{i-1}^{(k)}), \quad i = 2, 3, \dots, m.$$

步 3, 由式 (3.113) 计算 $\mathbf{y}_s^{(k)} = (\alpha_1^{(k)}, \dots, \alpha_{s-1}^{(k)}, \alpha_{s+1}^{(k)}, \dots, \alpha_m^{(k)})^T$ 及

$$\alpha_s^{(k)} = 1 - \sum_{i=1, i \neq s}^m \alpha_i^{(k)}.$$

步 4, 计算 $\mathbf{x}^{(k)} = \alpha_1^{(k)} \mathbf{x}_1^{(k)} + \alpha_2^{(k)} \mathbf{x}_2^{(k)} + \cdots + \alpha_m^{(k)} \mathbf{x}_m^{(k)}$.

步 5, 若 $\|\mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}\|/\|\mathbf{b}\| \leq \varepsilon$, 则停算, 输出 $\mathbf{x}^{(k)}$ 作为方程组的近似解.

步 6, 置 $k := k + 1$, 转步 2.

注 3.5 记 $\mathbf{r}_i^{(k)} = \mathbf{b} - \mathbf{A}\mathbf{x}_i^{(k)}$, 由式 (3.117) 可得: $\|\mathbf{r}^{(k)}\|_2$ 是坐标原点与“由点 $\mathbf{r}^{(k-1)}, \mathbf{r}_2^{(k)}, \dots, \mathbf{r}_m^{(k)}$ ”确定的超平面 π_k 之间的最短距离. 因此, $\|\mathbf{r}^{(k)}\|_2 < \|\mathbf{r}^{(k-1)}\|_2$ 等价于 π_k 不与“以 $\mathbf{r}^{(k-1)}$ 为法向量”的超平面

$$S_k = \{\mathbf{r}(\mathbf{z}) : \mathbf{r}^{(k-1)} \perp (\mathbf{r}(\mathbf{z}) - \mathbf{r}^{(k-1)}), \mathbf{z} \in \mathbb{C}^n\}$$

重合, 即存在 $i_0 \in \{2, 3, \dots, m\}$, 使 $\mathbf{r}_{i_0}^{(k)} \notin S_k$, 也就是

$$\mathbf{x}_{i_0}^{(k)} \notin T_k = \{\mathbf{z} : \mathbf{r}^{(k-1)} \perp (\mathbf{r}(\mathbf{z}) - \mathbf{r}^{(k-1)}), \mathbf{z} \in \mathbb{C}^n\}.$$

若用空间 \mathbb{C}^n 中的“体积”定义集合的测度, 那么超平面集合 T_k 的测度值为零. 因此, $\mathbf{x}_i^{(k)} (i = 2, 3, \dots, m)$ 几乎不同时落入集合 T_k , 从而可得: $\|\mathbf{r}^{(k)}\|_2 < \|\mathbf{r}^{(k-1)}\|_2$ 几乎对所有的 k 成立.

一方面, 当式 (3.120) 收敛时, 算法 3.9 必定收敛. 进一步, 若第 k 次校正过程的整体缩减系数

$$\alpha^{(k)} = \frac{\|\mathbf{r}^{(k)}\|_2}{\min_{1 \leq i \leq m} \|\mathbf{r}_i^{(k)}\|_2} < 1,$$

则算法 3.9 比式 (3.120) 收敛得快; 另一方面, 当式 (3.120) 不收敛时, 由于 $\|\mathbf{r}^{(i)}\|_2 < \|\mathbf{r}^{(i-1)}\|_2$ 几乎对所有的 i 成立, 所以算法 3.9 也能收敛.

例 3.3 设

$$\mathbf{A} = \begin{bmatrix} 1.2 & 0.3 & 0.4 \\ 0.4 & 1.2 & 0.3 \\ 0.3 & 0.4 & 1.2 \end{bmatrix}, \mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}.$$

取 $\mathbf{x}^{(0)} = (1, 2, 0)^T$, 式 (3.120) 采用 Richardson 迭代格式 (3.26), 其中

$$\mathbf{B} = \mathbf{I} - \mathbf{A} = - \begin{bmatrix} 0.2 & 0.3 & 0.4 \\ 0.4 & 0.2 & 0.3 \\ 0.3 & 0.4 & 0.2 \end{bmatrix}.$$

取终止准则值为

$$\frac{\|\mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}\|_2}{\|\mathbf{b}\|_2} \leq 10^{-5}.$$

整体校正法中的 $m = 2, s = 2$. 由 $\rho(\mathbf{B}) \leq \|\mathbf{B}\|_1 = 0.9$ 知格式 (3.4) 收敛, 而算法 3.9 收敛较快, 计算结果如表 3.3 所示.

方程组 $\mathbf{A}\mathbf{x} = \mathbf{b}$ 的精确解为 $\mathbf{x}^* = \left(\frac{1}{1.9}, \frac{1}{1.9}, \frac{1}{1.9}\right)^T$.

表 3.3 整体校正取 $m = 2, k = 2$

迭代法	迭代次数	CPU 时间	相对残差
Richardson 迭代	109	0.0004	9.2614e-06
算法 3.9	11	0.0004	2.9541e-06

例 3.4 设

$$A = \begin{bmatrix} 4 & -1 & & & \\ -1 & 4 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 4 & -1 \\ & & & -1 & 4 \end{bmatrix} \in \mathbb{R}^{1023 \times 1023}, \quad b = \begin{bmatrix} 3 \\ 2 \\ \vdots \\ 2 \\ 3 \end{bmatrix}.$$

取 $x^{(0)} = 0$, 式 (3.120) 仍采用 Richardson 迭代格式 (3.26), 其中

$$B = I - A = \begin{bmatrix} -3 & 1 & & & \\ 1 & -3 & 1 & & \\ & \ddots & \ddots & \ddots & \\ & & 1 & -3 & 1 \\ & & & 1 & -3 \end{bmatrix}.$$

取终止准则值为

$$\frac{\|b - Ax^{(k)}\|_2}{\|b\|_2} \leq 10^{-10}.$$

整体校正法中的 $m = 10, s = 10$. 由 $\rho(B) > 1$ 知格式 (3.4) 发散, 而算法 3.9 收敛. 计算结果如表 3.4 所示.

表 3.4 整体校正取 $m = 10, k = 10$

迭代法	迭代次数	CPU 时间	相对残差
Richardson 迭代	1000	0.4950	NaN
算法 3.9	5	0.1458	2.7583e-12

3.5.3 基于矩阵特征值的外推方法

考虑迭代格式 (3.4), 当 $\rho(B) < 1$, 但 $\rho(B) \approx 1$ 时, 格式收敛缓慢. 下面使用外推技术, 改善格式 (3.4) 的收敛性态.

设迭代矩阵 B 的 n 个实特征值满足 $1 > |\lambda_1| > |\lambda_2| \geq \cdots \geq |\lambda_n|$, 对应的特征向量分别为 z_1, z_2, \cdots, z_n , 且它们线性无关. 分解

$$x^{(1)} - x^{(0)} = c_1 z_1 + c_2 z_2 + \cdots + c_n z_n,$$

$$\mathbf{x}^{(i+1)} - \mathbf{x}^{(i)} = \mathbf{B}^i(\mathbf{x}^{(1)} - \mathbf{x}^{(0)}) = c_1 \lambda_1^i \mathbf{z}_1 + c_2 \lambda_2^i \mathbf{z}_2 + \cdots + c_n \lambda_n^i \mathbf{z}_n,$$

则有

$$\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} = c_1 \lambda_1^k \mathbf{z}_1 + O(\lambda_2^k).$$

由于 $\rho(\mathbf{B}) < 1$, 所以 $\lim_{m \rightarrow \infty} \mathbf{x}^{(m)} = \mathbf{x}^*$, 且有

$$\begin{aligned} \mathbf{x}^* - \mathbf{x}^{(k)} &= \lim_{m \rightarrow \infty} (\mathbf{x}^{(m)} - \mathbf{x}^{(k)}) = \lim_{m \rightarrow \infty} \sum_{i=k}^{m-1} (\mathbf{x}^{(i+1)} - \mathbf{x}^{(i)}) \\ &= \lim_{m \rightarrow \infty} \sum_{i=k}^{m-1} (c_1 \lambda_1^i \mathbf{z}_1 + c_2 \lambda_2^i \mathbf{z}_2 + \cdots + c_n \lambda_n^i \mathbf{z}_n) \\ &= \frac{\lambda_1^k}{1 - \lambda_1} c_1 \mathbf{z}_1 + \frac{\lambda_2^k}{1 - \lambda_2} c_2 \mathbf{z}_2 + \cdots + \frac{\lambda_n^k}{1 - \lambda_n} c_n \mathbf{z}_n. \end{aligned} \quad (3.121)$$

从而可得

$$\mathbf{x}^* = \mathbf{x}^{(k)} + \frac{1}{1 - \lambda_1} (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) + O(\lambda_2^k).$$

记

$$\tilde{\mathbf{x}}^{(k+1)} = \mathbf{x}^{(k)} + \frac{1}{1 - \lambda_1} (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}),$$

则有 $\tilde{\mathbf{x}}^{(k+1)} - \mathbf{x}^* = O(\lambda_2^k)$. 再由式 (3.121), 得 $\mathbf{x}^{(k+1)} - \mathbf{x}^* = O(\lambda_1^{k+1})$.

当 $|\lambda_2| \ll |\lambda_1| < 1$ 且 $|\lambda_1| \approx 1$ 时, 有

$$\frac{\lambda_2^k}{\lambda_1^{k+1}} = \left(\frac{\lambda_2}{\lambda_1} \right)^k \frac{1}{\lambda_1} \rightarrow 0 \quad (k \rightarrow \infty).$$

因此, $\tilde{\mathbf{x}}^{(k+1)}$ 比 $\mathbf{x}^{(k+1)}$ 更接近于 \mathbf{x}^* .

将外推技术施加于迭代格式 (3.4), 即可得到基于迭代矩阵特征值的外推格式.

算法 3.10 (基于特征值的外推法)

给定整数 $m \geq 1$, 初始值 $\mathbf{x}_1^{(0)}$ 及最大迭代次数 N .

while $i \leq N$

for $k = 1 : m$

$$\mathbf{x}_i^{(k)} = \mathbf{B} \mathbf{x}_i^{(k-1)} + \mathbf{f};$$

end

$$\mathbf{x}_{i+1}^{(0)} = \mathbf{x}_i^{(m-1)} + \frac{1}{1 - \lambda_1} (\mathbf{x}_i^{(m)} - \mathbf{x}_i^{(m-1)});$$

$$i = i + 1;$$

end

易知, 当 $|\lambda_2| \ll |\lambda_1| < 1$ 且 $|\lambda_1| \approx 1$ 时, 算法 3.10 比迭代格式 (3.4) 收敛得快.

例 3.5 在例 3.2 中, 用 Jacobi 迭代法收敛的速度相当缓慢. 现在对 Jacobi 迭代法进行加速. 先计算 Jacobi 迭代矩阵 B_J 的特征值. 注意到 B_J 可以表示为

$$B_J = D^{-1}(D - A) = \frac{1}{4}(2I - S), \quad (3.122)$$

式中: I 为 n 阶单位矩阵.

$$S = \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & & \ddots & \ddots & \ddots \\ & & & -1 & 2 & -1 \\ & & & & -1 & 2 \end{bmatrix}. \quad (3.123)$$

S 的第 k 个特征值为 $4\sin^2\left(\frac{k\pi}{2(n+1)}\right)$, $k = 1, 2, \dots, n$. 故 B_J 的第 k 个特征值为

$$\lambda_k = \frac{1}{4}\left[2 - 4\sin^2\left(\frac{k\pi}{2(n+1)}\right)\right] = \frac{1}{2}\cos\frac{k\pi}{n+1}, \quad k = 1, 2, \dots, n. \quad (3.124)$$

由式 (3.124) 可知 B_J 模最大特征值为 $\lambda_1 = \frac{1}{2}\cos\frac{\pi}{n+1}$. 取 $n = 2^{14} - 1 = 16383$, 初始向量为零向量, 容许误差为 10^{-10} . 利用算法 3.10, 得到计算结果如表 3.5 所示.

表 3.5 Jacobi 迭代法和特征值外推法的数值比较

迭代格式	迭代次数 (k)	CPU 时间	相对残差
Jacobi 迭代法	34	0.0158	5.8182e-11
特征值外推法	4	0.0076	1.1501e-12

3.5.4 Chebyshev 加速方法

考虑迭代格式 (3.4). 设 $\rho(B) < 1$, 由初始向量 $\mathbf{x}^{(0)}$ 出发, 经过 k 次迭代得到 $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(k)}$ ($k \geq 1$). 选择数组 $\alpha_0^{(k)}, \alpha_1^{(k)}, \dots, \alpha_k^{(k)}$, 使得

$$\mathbf{u}^{(k)} = \sum_{i=0}^k \alpha_i^{(k)} \mathbf{x}^{(i)}, \quad \sum_{i=0}^k \alpha_i^{(k)} = 1. \quad (3.125)$$

设 $\mathbf{x}^* = B\mathbf{x}^* + \mathbf{f}$, 则

$$\mathbf{u}^{(k)} - \mathbf{x}^* = \sum_{i=0}^k \alpha_i^{(k)} (\mathbf{x}^{(i)} - \mathbf{x}^*) = \sum_{i=0}^k \alpha_i^{(k)} B^i (\mathbf{x}^{(0)} - \mathbf{x}^*).$$

定义 k 次多项式

$$p_k(t) = \sum_{i=0}^k \alpha_i^{(k)} t^i, \quad p_k(1) = 1,$$

即有

$$\mathbf{u}^{(k)} - \mathbf{x}^* = p_k(B)(\mathbf{x}^{(0)} - \mathbf{x}^*).$$

根据上式, 有

$$\|u^{(k)} - x^*\|_2 \leq \|p_k(B)\|_2 \cdot \|x^{(0)} - x^*\|_2.$$

对所有的 $\|p_k(B)\|_2$ 取下确界就是 $\rho(p_k(B))$. 很自然地, 要求选取多项式 p_k 使得 $\rho(p_k(B))$ 达到最小. 设 S 是复平面上包含 B 的所有特征值的集合, 则

$$\rho(p_k(B)) \leq \max_{\lambda \in \lambda(B)} |p_k(\lambda)| \leq \max_{\lambda \in S} |p_k(\lambda)|.$$

设迭代矩阵 B 的 n 个特征值 $\lambda_1, \lambda_2, \dots, \lambda_n$ 均为实数, 都在区间 $[a, b]$ 上, 并且 $1 \notin [a, b]$. 这里只考虑 $a < b < 1$ 的情形. 为了使 $\rho(p_k(B))$ 达到极小, 只要考虑极小极大问题

$$\min_{p_k \in \mathcal{P}_k^{(1)}} \left\{ \max_{\lambda \in [a, b]} |p_k(\lambda)| \right\}, \quad (3.126)$$

式中: $\mathcal{P}_k^{(1)}$ 为满足 $p_k(1) = 1$ 所有次数不超过 k 的实系数多项式的集合. 这是一个在多项式空间有归一化条件的最佳一致逼近问题, 它与 Chebyshev 多项式有关. k 次 Chebyshev 多项式 $C_k(t)$ 是极小极大问题

$$\min_{p_k \in \mathcal{P}_k^{(1)}} \left\{ \max_{t \in [-1, 1]} |p_k(t)| \right\} \quad (3.127)$$

的解. k 次 Chebyshev 多项式的定义为

$$C_k(t) = \begin{cases} \cos(k \cdot \arccos t), & |t| \leq 1, \\ \cosh(k \cdot \operatorname{arccosh} t), & |t| > 1, \end{cases} \quad (3.128)$$

它满足递推关系

$$\begin{cases} C_0(t) = 1, \quad C_1(t) = t, \\ C_{k+1}(t) = 2tC_k(t) - C_{k-1}(t), \quad k \geq 1. \end{cases} \quad (3.129)$$

下面的引理给出了在归一化条件下极小极大问题 (3.127) 的最优解.

引理 3.5 设 $\beta \notin [-1, 1]$. 若 $p_k \in \mathcal{P}_k$, $p_k(\beta) = 1$, 则

$$\|p_k\| \equiv \max_{-1 \leq t \leq 1} |p_k(t)| \geq \frac{1}{C_k(\beta)}.$$

特别地, 当取 $p_k(t) = \frac{C_k(t)}{C_k(\beta)}$ 时, $\|p_k\| = \frac{1}{C_k(\beta)}$.

证明 令 $\alpha = \frac{1}{C_k(\beta)}$. 不难求得 $C_k(t)$ 的极值点为 $t_i = \cos\left(\frac{i\pi}{k}\right)$, 满足 $C_k(t_i) = \cos(i\pi) = (-1)^i$, $i = 0, 1, \dots, k$. 用反证法. 若 $\|p_k\| < \alpha$, 则有

$$\operatorname{sign}(\alpha)(-1)^i(\alpha C_k(t_i) - p_k(t_i)) \geq |\alpha| - \|p_k\| > 0,$$

这表明多项式 $\alpha C_k(t) - p_k(t)$ 在 $t_i (i = 0, 1, \dots, k)$ 处有交替的正负号, 所以它在 $(-1, 1)$ 中有 k 个零点. 但 $\beta \notin [-1, 1]$ 也是 $\alpha C_k(t) - p_k(t)$ 的零点, 故 k 次多项式 $\alpha C_k(t) - p_k(t)$ 有 $k+1$ 个零点, 矛盾. 因此, 必有 $\|p_k\| \geq \alpha$. 另外, 显然当 $p_k(t) = C_k(t)/C_k(\beta)$ 时, $\|p_k\| = \alpha$. 证毕. \square

下面的引理给出了极小极大问题 (3.126) 的最优解.

引理 3.6 设 $a < b < 1$. 若 $p_k \in \mathcal{P}_k$ 且满足 $p_k(1) = 1$, 则

$$\|p_k\| \equiv \max_{a \leq t \leq b} |p_k(t)| \geq \frac{1}{C_k(w(1))},$$

式中:

$$w(t) = \frac{2t - a - b}{b - a}.$$

特别地, 当取 $p_k(t) = \frac{C_k(w(t))}{C_k(w(1))}$ 时, $\|p_k\| = \frac{1}{C_k(w(1))}$.

证明 记 $\beta = w(1) = \frac{2 - a - b}{b - a} > 1$. 注意到 $w(t) : [a, b] \rightarrow [-1, 1]$, 取 $q_k \in \mathcal{P}_k$, 使得 $p_k(t) = q_k(w(t))$. 显然有 $1 = p_k(1) = q_k(w(1)) = q_k(\beta)$. 根据引理 3.5, 有

$$\|q_k\| = \max_{-1 \leq w(t) \leq 1} |q_k(w(t))| \geq \frac{1}{C_k(\beta)},$$

即

$$\|p_k\| = \max_{a \leq t \leq b} |p_k(t)| \geq \frac{1}{C_k(\beta)} = \frac{1}{C_k(w(1))}.$$

此外, 显然当 $p_k(t) = \frac{C_k(w(t))}{C_k(w(1))}$ 时, $\|p_k\| = \frac{1}{C_k(w(1))}$. 证毕. \square

注 3.6 若 $1 < a < b$, 只需作变换 $w(t) : [a, b] \rightarrow [-1, 1]$,

$$w(t) = \frac{2t - a - b}{a - b},$$

则有

$$\beta = w(1) = \frac{2 - a - b}{a - b} = 1 + 2 \frac{1 - a}{a - b} > 1,$$

从而亦可推导出, 当 $p_k(t) = \frac{C_k(w(t))}{C_k(w(1))}$ 时, $\|p_k\| = \frac{1}{C_k(w(1))}$.

下面的引理给出了最优解 p_k 的递推关系.

引理 3.7 $p_k(t) = C_k(w(t))/C_k(w(1))$ 满足如下的递推关系:

$$\begin{cases} p_0(t) = 1, & p_1(t) = w(t)/w(1) = (2t - a - b)/(2 - a - b), \\ p_k(t) = \rho_k p_1(t) p_{k-1}(t) + (1 - \rho_k) p_{k-2}(t), & k \geq 2, \end{cases} \quad (3.130)$$

式中: ρ_k 满足的递推关系为

$$\rho_1 = 2, \quad \rho_k = (1 - \alpha \rho_{k-1})^{-1}, \quad \alpha = (2w(1))^{-2}, \quad k \geq 2. \quad (3.131)$$

证明 令 $\beta_k = C_k(w(1))$. 根据 C_k 的递推关系, 得

$$\begin{aligned} p_k(t) &= \beta_k^{-1} C_k(w(t)) = \beta_k^{-1} [2w(t)C_{k-1}(w(t)) - C_{k-2}(w(t))] \\ &= 2\beta_k^{-1} \beta_{k-1} w(1) p_1(t) p_{k-1}(t) - \beta_k^{-1} \beta_{k-2} p_{k-2}(t). \end{aligned} \quad (3.132)$$

定义

$$\rho_k = 2\beta_k^{-1} \beta_{k-1} w(1) = \alpha^{-1/2} \beta_k^{-1} \beta_{k-1}. \quad (3.133)$$

则由 C_k 的递推过程可以推出 β_k 的递推关系为

$$\beta_k = C_k(w(1)) = 2w(1)C_{k-1}(w(1)) - C_{k-2}(w(1)) = \alpha^{-1/2} \beta_{k-1} - \beta_{k-2}.$$

从而有

$$\beta_k^{-1} \beta_{k-2} = \beta_k^{-1} (\alpha^{-1/2} \beta_{k-1} - \beta_k) = \rho_k - 1. \quad (3.134)$$

将式 (3.133) 和式 (3.134) 代入式 (3.132) 即可得到 p_k 的递推公式 (3.130). ρ_k 的递推关系亦可从 β_k 的递推过程中得到, 即由式 (3.133), 得

$$\begin{aligned} \rho_k &= \alpha^{-1/2} \beta_{k-1} \beta_k^{-1} = \alpha^{-1/2} \beta_{k-1} (\alpha^{-1/2} \beta_{k-1} - \beta_{k-2})^{-1} \\ &= (1 - \alpha^{1/2} \beta_{k-1}^{-1} \beta_{k-2})^{-1} = (1 - \alpha \rho_{k-1})^{-1}. \end{aligned}$$

证毕. □

引理 3.7 的意义在于表明了由 Chebyshev 加速方法得到的序列 $\{u^{(k)}\}$ 也有递推公式, 从而在加速过程中没有必要真正地引入 Chebyshev 多项式. 有下面的定理.

定理 3.36 设 $\{u^{(k)}\}$ 是由 Chebyshev 加速方法得到的序列, 则对于任意的 $u^{(0)}$, 有如下的递推公式

$$\begin{cases} u^{(1)} = (1 - \gamma)u^{(0)} + \gamma(Bu^{(0)} + f), \quad \gamma = 2/(2 - a - b), \\ u^{(k)} = (1 - \rho_k)u^{(k-2)} + \rho_k[(1 - \gamma)u^{(k-1)} + \gamma(Bu^{(k-1)} + f)], \quad k \geq 2, \end{cases} \quad (3.135)$$

式中: 序列 $\{\rho_k\}$ 由式 (3.131) 所定义.

证明 由于

$$u^{(k)} = \sum_{i=0}^k \alpha_i^{(k)} x^{(i)}, \quad p_k(t) = \sum_{i=0}^k \alpha_i^{(k)} t^i, \quad \sum_{i=0}^k \alpha_i^{(k)} = 1,$$

故 $u^{(0)} = x^{(0)}$, $u^{(1)}$ 由式 (3.135) 的第 1 式给出. 设 $u^* = Bu^* + f$, 则

$$\begin{aligned} u^{(k)} - u^* &= p_k(B)(u^{(0)} - u^*) \\ &= [\rho_k p_1(B) p_{k-1}(B) + (1 - \rho_k) p_{k-2}(B)](u^{(0)} - u^*) \\ &= \rho_k p_1(B)(u^{(k-1)} - u^*) + (1 - \rho_k)(u^{(k-2)} - u^*), \end{aligned}$$

即

$$\mathbf{u}^{(k)} = \rho_k p_1(\mathbf{B}) \mathbf{u}^{(k-1)} + (1 - \rho_k) \mathbf{u}^{(k-2)} + \rho_k (\mathbf{I} - p_1(\mathbf{B})) \mathbf{x}^*, \quad k \geq 2. \quad (3.136)$$

注意到

$$p_1(\mathbf{B}) = (1 - \gamma) \mathbf{I} + \gamma \mathbf{B}, \quad (\mathbf{I} - \mathbf{B}) \mathbf{x}^* = \mathbf{f},$$

代入式 (3.136) 即得要证的递推式. 证毕. \square

注 3.7 式 (3.135) 中的迭代格式表明, 对于 $\mathbf{u}^{(k-1)}$, 用原来的迭代格式迭代一次产生 $\mathbf{B}\mathbf{u}^{(k-1)} + \mathbf{f}$, 它与 $\mathbf{u}^{(k-1)}$ 组合产生 $(1 - \gamma) \mathbf{u}^{(k-1)} + \gamma (\mathbf{B}\mathbf{u}^{(k-1)} + \mathbf{f})$, 其中组合系数 γ 由 a, b 决定, 这种组合就是外推格式. 这样外推之后的结果再与前一步迭代向量 $\mathbf{u}^{(k-2)}$ 做组合, 组合系数 ρ_k 由式 (3.131) 的递推公式确定. 可以期望, Chebyshev 加速应该比一般的外推加速具有更快的收敛速度.

Chebyshev 加速迭代格式的速度由它的迭代矩阵的谱半径确定. 有

$$\rho(p_k(\mathbf{B})) \leq \max_{a \leq t \leq b} |p_k(\lambda)| = \frac{1}{C_k(w(1))}.$$

经过简单的计算 (定理 4.2 (2)) 可知

$$\frac{1}{C_k(w(1))} = \frac{2}{r^k + r^{-k}}, \quad r = t + \sqrt{t^2 - 1}, \quad t = w(1) = \frac{2 - a - b}{b - a} > 1.$$

外推格式和 Chebyshev 加速迭代格式都依赖于某个简单的迭代格式 (3.4). 当迭代格式 (3.4) 容易并行实现时, 相应的外推格式和 Chebyshev 加速迭代格式也都可以并行实现.

例 3.6 考虑用 Chebyshev 加速方法对例 3.2 中的 Jacobi 迭代法进行加速. 现在对 Jacobi 迭代法进行加速. 根据例 3.5, 可知 Jacobi 迭代矩阵 \mathbf{B}_J 的 n 个特征值为

$$\lambda_k = \frac{1}{2} \cos \frac{k\pi}{n+1}, \quad k = 1, 2, \dots, n.$$

因此, \mathbf{B}_J 所有特征值都在 $a = \lambda_n$ 和 $b = \lambda_1$ 之间. 取 $n = 2^{14} - 1 = 16383$, 初始向量为零向量, 容许误差为 10^{-10} . 由于 $\lambda_n = -\lambda_1$, 故外推参数 $\gamma = 2/(2 - a - b) = 1$, 因此外推法 (3.108) 对此例不起作用. 得到计算结果如表 3.6 所示.

表 3.6 Jacobi 迭代法和 Chebyshev 加速迭代的数值比较

迭代格式	迭代次数 (k)	CPU 时间	相对残差
Jacobi 迭代法	34	0.0130	5.8182e-11
Chebyshev 加速	21	0.0089	6.0411e-11

3.6 块三对角方程组的迭代解法

本节讨论 A 为块三对角矩阵时, 方程组 $Ax = b$ 的迭代解法. 设

$$A = \begin{bmatrix} B_1 & C_1 & & & \\ A_2 & B_2 & C_2 & & \\ & \ddots & \ddots & \ddots & \\ & & A_{m-1} & B_{m-1} & C_{m-1} \\ & & & A_m & B_m \end{bmatrix}, \quad (3.137)$$

式中: B_i 为 n_i 阶矩阵 ($n_1 + n_2 + \cdots + n_m = n$).

3.6.1 PE(α) 方法

对矩阵 A 进行块三角分解 (假定这种分解存在), 即

$$A = \begin{bmatrix} D_1 & & & & \\ A_2 & D_2 & & & \\ & \ddots & \ddots & & \\ & & A_m & D_m \end{bmatrix} \begin{bmatrix} I_1 & U_1 & & & \\ & \ddots & \ddots & & \\ & & I_{m-1} & U_{m-1} & \\ & & & I_m \end{bmatrix}, \quad (3.138)$$

式中: I_i 为 n_i 阶单位矩阵.

比较式 (3.137) 和式 (3.138), 得

$$D_1 = B_1,$$

$$U_i = D_i^{-1} C_i \quad (i = 1, 2, \cdots, m-1),$$

$$D_i = B_i - A_i U_{i-1} = B_i - A_i D_{i-1}^{-1} C_{i-1} \quad (i = 2, 3, \cdots, m).$$

由此, 在一定条件下, 有

$$\begin{aligned} D_i^{-1} &= (B_i - A_i D_{i-1}^{-1} C_{i-1})^{-1} = (I_i - B_i^{-1} A_i D_{i-1}^{-1} C_{i-1})^{-1} B_i^{-1} \\ &= [I_i + B_i^{-1} A_i D_{i-1}^{-1} C_{i-1} + (B_i^{-1} A_i D_{i-1}^{-1} C_{i-1})^2 + \cdots] B_i^{-1} \\ &\approx (I_i + B_i^{-1} A_i D_{i-1}^{-1} C_{i-1}) B_i^{-1}. \end{aligned}$$

若将右端的 $B_i^{-1} A_i D_{i-1}^{-1} C_{i-1}$ 也省略, 则有 $D_i^{-1} \approx B_i^{-1}$. 于是, 用 B_{i-1}^{-1} 代替 D_{i-1}^{-1} 时, 上式可写为 $D_i^{-1} \approx (I_i + B_i^{-1} A_i B_{i-1}^{-1} C_{i-1}) B_i^{-1}$, 或

$$D_i \approx B_i (I_i + B_i^{-1} A_i B_{i-1}^{-1} C_{i-1})^{-1} = S_i, \quad (3.139)$$

这样得到 A 的一个近似块三角分解

$$A \approx \begin{bmatrix} S_1 & & & & \\ A_2 & S_2 & & & \\ & \ddots & \ddots & & \\ & & A_m & S_m \end{bmatrix} \begin{bmatrix} I_1 & T_1 & & & \\ & \ddots & \ddots & & \\ & & I_{m-1} & T_{k-1} & \\ & & & I_m \end{bmatrix}, \quad (3.140)$$

式中:

$$\begin{cases} S_1 = B_1, \\ T_i = S_i^{-1}C_i \quad (i = 1, 2, \dots, m-1), \\ S_i = B_i(I_i + B_i^{-1}A_iB_{i-1}^{-1}C_{i-1})^{-1} \quad (i = 2, 3, \dots, m). \end{cases} \quad (3.141)$$

若用 L 和 U 分别表示式 (3.140) 右端的第 1 个和第 2 个矩阵, 则有 $A \approx LU$. 令

$$M = LU = \begin{bmatrix} S_1 & C_1 & & & \\ A_2 & A_2T_1 + S_2 & C_2 & & \\ & \ddots & \ddots & \ddots & \\ & & A_{m-1} & A_{m-1}T_{m-2} + S_{m-1} & C_{m-1} \\ & & & A_m & A_mT_{m-1} + S_m \end{bmatrix},$$

$$N = M - A = \text{diag}(O, N_2, \dots, N_m), \quad (3.142)$$

式中: $N_i = A_iT_{i-1} + S_i - B_i$ ($i = 2, 3, \dots, m$). 那么 $Ax = b$ 等价于 $Mx = Nx + b$. 建立迭代格式

$$Mx^{(k+1)} = Nx^{(k)} + b, \quad k = 0, 1, 2, \dots, \quad (3.143)$$

称式 (3.143) 为拟消去 (Pseudo-Elimination, PE) 迭代方法. 使用式 (3.143) 时, 对于每个 k , 求 $x^{(k+1)}$ 可转化为下面两个块三角方程组的求解问题:

$$Lz^{(k+1)} = Nx^{(k)} + b, \quad Ux^{(k+1)} = z^{(k+1)}.$$

若用 $b_i, x_i^{(k)}, x_i^{(k+1)}, z_i^{(k+1)}$ 分别表示向量 $b, x^{(k)}, x^{(k+1)}, z^{(k+1)}$ 的第 i 个子向量, 则分组计算的递推公式为

$$z_1^{(k+1)} = S_1^{-1}b_1, \quad z_i^{(k+1)} = S_i^{-1}(N_i x_i^{(k)} + b_i - A_i z_{i-1}^{(k+1)}), \quad i = 2, 3, \dots, m,$$

$$x_m^{(k+1)} = z_m^{(k+1)}, \quad x_i^{(k+1)} = z_i^{(k+1)} - T_i x_{i+1}^{(k+1)}, \quad i = m-1, \dots, 2, 1.$$

下面讨论 PE 方法的收敛性.

定理 3.37 设 A 为 Hermite 矩阵, 则 PE 方法中的矩阵 S_i, M 和 N 均为 Hermite 矩阵.

证明 由 $A^H = A$ 知 $B_i^H = B_i, A_i^H = C_{i-1}$, 故有

$$\begin{aligned} S_i^H &= [(I_i + B_i^{-1}A_iB_{i-1}^{-1}C_{i-1})^{-1}]^H B_i^H = (I_i + A_iB_{i-1}^{-1}C_{i-1}B_i^{-1})^{-1}B_i \\ &= [B_i(I_i + B_i^{-1}A_iB_{i-1}^{-1}C_{i-1})B_i^{-1}]^{-1}B_i = B_i(I_i + B_i^{-1}A_iB_{i-1}^{-1}C_{i-1})^{-1} \\ &= S_i \end{aligned}$$

及

$$(A_iT_{i-1})^H = (A_iS_{i-1}^{-1}C_{i-1})^H = C_{i-1}^H S_{i-1}^{-1} A_i^H = A_i S_{i-1}^{-1} C_{i-1} = A_i T_{i-1}.$$

因此 $M^H = M$, 从而 $N^H = (M - A)^H = M^H - A^H = N$. 证毕. \square

引理 3.8 设 A 的特征值 $\lambda(A)$ 为非负实数, B 为 Hermite 正定矩阵, $C = AB$ (或 $C = BA$) 为 Hermite 矩阵, 则 C 为半正定矩阵. 进一步, 若 $\lambda(A)$ 均为正数, 则 C 为正定矩阵.

证明 设 $\lambda(A) \geq 0$, 由于 A 相似于

$$B^{-\frac{1}{2}}AB^{\frac{1}{2}} = B^{-\frac{1}{2}}ABB^{-\frac{1}{2}} = B^{-\frac{1}{2}}CB^{-\frac{1}{2}},$$

所以 $B^{-\frac{1}{2}}CB^{-\frac{1}{2}}$ 的特征值也是非负实数, 从而 C 的特征值是非负实数. 故 C 为半正定矩阵. 同理, 可证 $C = BA$ 及 $\lambda(A) > 0$ 的情形. 证毕. \square

引理 3.9 设 $A = M - N$, 且 A 和 M 都可逆, 则对 $M^{-1}N$ 的任意特征值 μ , 存在 $A^{-1}N$ 的某个特征值 λ , 使得

$$\mu = \frac{\lambda}{1 + \lambda}. \quad (3.144)$$

证明 设 $M^{-1}N$ 的一个特征值为 μ , 对应的特征向量为 $z \neq 0$, 则有 $(M^{-1}N)z = \mu z$. 因为

$$M^{-1}N = (A + N)^{-1}N = (I + A^{-1}N)^{-1}(A^{-1}N),$$

所以

$$(I + A^{-1}N)^{-1}(A^{-1}N)z = \mu z,$$

即

$$(A^{-1}N)z = \mu(I + A^{-1}N)z.$$

若 $\mu = 1$, 则上式成为 $z = 0$, 这与 z 是特征向量矛盾. 故 $\mu \neq 1$, 从而有

$$(A^{-1}N)z = \frac{\mu}{1 - \mu}z,$$

这表明 $\frac{\mu}{1 - \mu}$ 是 $A^{-1}N$ 的一个特征值, 记作 λ , 即 $\lambda = \frac{\mu}{1 - \mu}$. 易见 $\lambda \neq -1$, 从而有 $\mu = \frac{\lambda}{1 + \lambda}$. 证毕. \square

定理 3.38 设 A 为 Hermite 正定矩阵, 则 PE 方法收敛.

证明 (1) 先证 PE 方法是适定的, 即算法中的矩阵 S_i ($i = 1, 2, \dots, m$) 均可逆. 由 A 正定知 B_i 均为正定矩阵, 从而 $A_i B_{i-1}^{-1} C_{i-1} = A_i B_{i-1}^{-1} A_i^H$ 半正定, 这是因为 $z \neq 0$ 时, 有

$$z^H A_i B_{i-1}^{-1} A_i^H z = (A_i^H z)^H B_{i-1}^{-1} (A_i^H z) \geq 0.$$

又由

$$B_i^{\frac{1}{2}}(B_i^{-1} A_i B_{i-1}^{-1} C_{i-1}) B_i^{-\frac{1}{2}} = B_i^{-\frac{1}{2}}(A_i B_{i-1}^{-1} C_{i-1}) B_i^{-\frac{1}{2}},$$

知 $B_i^{-1} A_i B_{i-1}^{-1} C_{i-1}$ 的特征值为非负实数, 从而 $I_i + B_i^{-1} A_i B_{i-1}^{-1} C_{i-1}$ 的特征值不为零 (实际上不小于 1). 因此 S_i 均可逆.

(2) 再证 PE 方法是收敛的, 即迭代矩阵 $M^{-1}N$ 的谱半径 $\rho(M^{-1}N) < 1$. 由于

$$\begin{aligned} N_i &= A_i T_{i-1} + S_i - B_i = A_i S_{i-1}^{-1} C_{i-1} + S_i - B_i \\ &= A_i (I_{i-1} + B_{i-1}^{-1} A_{i-1} B_{i-2}^{-1} C_{i-2}) B_{i-1}^{-1} C_{i-1} \\ &\quad + B_i (I_i + B_i^{-1} A_i B_{i-1}^{-1} C_{i-1})^{-1} - B_i \\ &\quad (i = 2, 3, \dots, m; A_1 = O, C_0 = O). \end{aligned}$$

利用恒等式

$$(I + G)^{-1} = I - G + G^2(I + G)^{-1},$$

改写上式第 2 项中的逆矩阵, 整理, 得

$$\begin{aligned} N_i &= A_i B_{i-1}^{-1} A_{i-1} B_{i-2}^{-1} C_{i-2} B_{i-1}^{-1} C_{i-1} \\ &\quad + A_i B_{i-1}^{-1} C_{i-1} B_i^{-1} A_i B_{i-1}^{-1} C_{i-1} (I + B_i^{-1} A_i B_{i-1}^{-1} C_{i-1})^{-1}. \end{aligned} \quad (3.145)$$

由 $A^H = C_{i-1}$ 知, 上式等号右端第 1 项为 Hermite 矩阵. 由定理 3.37 知 N_i 均为 Hermite 矩阵, 从而上式等号右端第 2 项也是 Hermite 矩阵. 由于

$$\begin{aligned} V_i &= A_i B_{i-1}^{-1} A_{i-1} B_{i-2}^{-1} C_{i-2} B_{i-1}^{-1} C_{i-1} \\ &= (A_i B_{i-1}^{-1} A_{i-1}) B_{i-2}^{-1} (A_i B_{i-1}^{-1} A_{i-1})^H, \\ W_i &= A_i B_{i-1}^{-1} C_{i-1} B_i^{-1} A_i B_{i-1}^{-1} C_{i-1} \\ &= (A_i B_{i-1}^{-1} A_i^H) B_i^{-1} (A_i B_{i-1}^{-1} A_i^H)^H, \\ Z_i &= B_i + A_i B_{i-1}^{-1} C_{i-1} = B_i + A_i B_{i-1}^{-1} A_i^H. \end{aligned}$$

所以 V_i 和 W_i 为半正定矩阵, Z_i 为正定矩阵. 再由

$$W_i Z_i^{-1} = Z_i^{\frac{1}{2}} (Z_i^{-\frac{1}{2}} W_i Z_i^{-\frac{1}{2}}) Z_i^{-\frac{1}{2}},$$

知 $W_i Z_i^{-1}$ 的特征值为非负实数. 式 (3.145) 等号右端的第 2 项是 Hermite 矩阵, 且可写为 $(W_i Z_i^{-1}) B_i$, 而 B_i 正定, 根据引理 3.8 可得 $(W_i Z_i^{-1}) B_i$ 为半正定矩阵. 于是有

$$\begin{aligned} N_i &= V_i + (W_i Z_i^{-1}) B_i \text{ 半正定} \\ \Rightarrow N &= \text{diag}(O, N_2, \dots, N_m) \text{ 半正定} \\ \Rightarrow A^{-1} N &= A^{-\frac{1}{2}} (A^{-\frac{1}{2}} N A^{-\frac{1}{2}}) A^{\frac{1}{2}} \text{ 的特征值为非负实数} \\ \Rightarrow 0 &\leq \frac{\lambda(A^{-1} N)}{1 + \lambda(A^{-1} N)} < 1. \end{aligned}$$

由式 (3.140) 知 M 可逆, 根据引理 3.9, 得 $\rho(M^{-1}N) < 1$. 证毕. □

在式 (3.139) 中, 改写 S_i 为

$$S_i = B_i (I_i + \alpha B_i^{-1} A_i B_{i-1}^{-1} C_{i-1})^{-1} \quad (i = 2, 3, \dots, m), \quad (3.146)$$

式中: α 为参数, 相应的迭代格式 (3.143) 称为 $PE(\alpha)$ 方法. 可以验证: $PE(0)$ 方法为对称块 Gauss-Seidel 方法 (简记为 SBGS 方法).

定理 3.39 设 A 为 Hermite 正定矩阵, 令 $G_i = B_i^{-1}A_iB_{i-1}^{-1}C_{i-1}$, $\delta = \lambda_{\min}(G_i)$, 若参数 α 满足

$$\begin{cases} \alpha \geqslant 0, & \text{当 } \delta \geqslant 1 \text{ 时,} \\ 0 \leqslant \alpha \leqslant \frac{1}{1-\delta}, & \text{当 } 0 \leqslant \delta < 1 \text{ 时,} \end{cases} \tag{3.147}$$

则 $PE(\alpha)$ 方法收敛.

在定理 3.39 中, G_i 的特征值为非负实数. 因此, 当 $\alpha \geqslant 0$ 时, $PE(\alpha)$ 方法收敛. 特别地, 当 $\alpha = 1$ 时, PE 方法收敛; 当 $\alpha = 0$ 时, SBGS 方法收敛.

例 3.7 设 A 形如式 (3.137), 其中

$$B_i = \begin{bmatrix} 4 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 4 \end{bmatrix}, \quad A_i = C_i = \begin{bmatrix} -1 & & \\ & -1 & \\ & & -1 \end{bmatrix}.$$

可以验证 A 为对称正定矩阵. 取 $b = (1, 2, \cdots, n)^T$, 用 $PE(\alpha)$ 方法求解线性方程组 $Ax = b$, 计算结果如表 3.7 所示 (终止准则为 10^{-12} , 时间单位为秒).

表 3.7 $PE(\alpha)$ 方法的数值结果

参数 α	$n = 6000$			$n = 12000$		
	迭代次数	计算时间	相对残差	迭代次数	计算时间	相对残差
0.5	22	0.8709	6.2931e-13	23	1.8432	3.3271e-13
1.0	15	0.6299	4.1516e-13	15	1.2642	7.6464e-13
1.4	8	0.3770	7.6054e-13	9	0.8406	4.6594e-14
1.5	7	0.3410	1.5256e-13	7	0.6851	2.8009e-13
1.6	9	0.4082	4.6582e-14	9	0.8277	8.5260e-14
2.0	17	0.6899	2.5846e-13	17	1.4063	4.7639e-13
SBGS	29	1.1470	8.4202e-13	30	2.4169	9.4989e-13

表 3.7 中最后一行是用 SBGS 方法求解该方程组时的迭代次数、计算时间和相对残差. 易见, 当参数 $\alpha = 1.5$ 时, $PE(\alpha)$ 方法的收敛速度较快.

3.6.2 二次 $PE(\alpha)$ 方法

在式 (3.139) 中, 改写 S_i 为

$$S_i = B_i(I_i + G_i + G_i^2)^{-1} \quad (i = 2, 3, \cdots, m), \tag{3.148}$$

式中: $G_i = B_i^{-1}A_iB_{i-1}^{-1}C_{i-1}$. 相应的迭代格式 (3.143) 称为二次 PE 方法.

定理 3.40 设 A 为 Hermite 矩阵, 则二次 PE 方法收敛.

证明 (1) 先证算法是适定的, 即所有的 S_i 是可逆的. 事实上, 由 A 正定知 B_i ($i = 1, 2, \dots, m$) 正定, 故 $S_1 = B_1$ 可逆. 而由 $A_i^H = C_{i-1}$ 知 $A_i B_{i-1}^{-1} C_{i-1} = A_i B_{i-1}^{-1} A_i^H$ 半正定. 再由

$$B_i^{-1} A_i B_{i-1}^{-1} C_{i-1} = B_i^{-\frac{1}{2}} \cdot B_i^{-\frac{1}{2}} (A_i B_{i-1}^{-1} C_{i-1}) B_i^{-\frac{1}{2}} B_i^{\frac{1}{2}}, \quad (3.149)$$

得 $B_i^{-1} A_i B_{i-1}^{-1} C_{i-1}$ 的特征值为非负实数, 从而

$$I_i + (B_i^{-1} A_i B_{i-1}^{-1} C_{i-1}) + (B_i^{-1} A_i B_{i-1}^{-1} C_{i-1})^2$$

的特征值为正数, 故该矩阵可逆. 根据式 (3.148), 得 S_i ($i = 1, 2, \dots, m$) 可逆.

(2) 再证算法是收敛的, 即迭代矩阵的谱半径 $\rho(M^{-1}N) < 1$. 由于

$$\begin{aligned} N_i &= A_i T_{i-1} + S_i - B_i = A_i S_{i-1}^{-1} C_{i-1} + S_i - B_i \\ &= A_i [I_{i-1} + (B_{i-1}^{-1} A_{i-1} B_{i-2}^{-1} C_{i-2}) + (B_{i-1}^{-1} A_{i-1} B_{i-2}^{-1} C_{i-2})^2] B_{i-1}^{-1} C_{i-1} \\ &\quad + B_i [I_i + (B_i^{-1} A_i B_{i-1}^{-1} C_{i-1}) + (B_i^{-1} A_i B_{i-1}^{-1} C_{i-1})^2]^{-1} - B_i. \end{aligned}$$

利用恒等式 $(I_i - G_i)(I_i + G_i + G_i^2) = I_i - G_i^3$, 得

$$(I_i + G_i + G_i^2)^{-1} = I_i - G_i + G_i^3(I_i + G_i + G_i^2)^{-1}.$$

代入上式并整理, 得

$$\begin{aligned} N_i &= A_i (B_{i-1}^{-1} A_{i-1} B_{i-2}^{-1} C_{i-2}) B_{i-1}^{-1} C_{i-1} \\ &\quad + A_i (B_{i-1}^{-1} A_{i-1} B_{i-2}^{-1} C_{i-2})^2 B_{i-1}^{-1} C_{i-1} + B_i G_i^3 (I_i + G_i + G_i^2)^{-1}. \end{aligned} \quad (3.150)$$

可以验证, 式 (3.150) 右端的前两项均为 Hermite 半正定矩阵. 由式 (3.149) 知, 矩阵 $G_i^3(I_i + G_i + G_i^2)^{-1}$ 的特征值为非负实数, 而 B_i 为 Hermite 正定矩阵, 且

$$\begin{aligned} B_i G_i^3(I_i + G_i + G_i^2)^{-1} &= B_i [(I_i + G_i + G_i^2)^{-1} + G_i - I_i] \\ &= (B_i^{-1} + G_i B_i^{-1} + G_i^2 B_i^{-1})^{-1} + B_i G_i - B_i \end{aligned}$$

为 Hermite 矩阵, 由引理 3.8 知式 (3.150) 右端的第 3 项为半正定矩阵. 因此, N_i 为半正定矩阵, 从而 N 亦为半正定矩阵. 再由 $A^{-1}N = A^{-\frac{1}{2}}(A^{-\frac{1}{2}}NA^{-\frac{1}{2}})A^{\frac{1}{2}}$ 可知 $\lambda(A^{-1}N)$ 为非负实数, 根据结论 (1) 可得 M 为可逆矩阵, 于是有

$$M^{-1}N = (N + A)^{-1}N = (A^{-1}N + I)^{-1}(A^{-1}N).$$

由此, 得

$$\lambda(M^{-1}N) = \frac{\lambda(A^{-1}N)}{1 + \lambda(A^{-1}N)}.$$

因此 $0 \leq \lambda(M^{-1}N) < 1$, 故迭代格式 (3.143) 收敛. 证毕. □

在式 (3.148) 中, 改写 S_i 为

$$S_i = B_i(I_i + G_i + \alpha G_i^2)^{-1} \quad (i = 2, 3, \dots, m), \quad (3.151)$$

式中: α, β 为参数. 相应的迭代格式 (3.143) 称为二次 $PE(\alpha)$ 方法. 易见, $PE(\alpha)$ 方法和二次 PE 方法都是二次 $PE(\alpha)$ 方法的特殊情形. 特别地, 二次 $PE(0)$ 方法即为 PE 方法.

定理 3.41 设 A 为 Hermite 正定矩阵, 若参数 α 满足式 (3.147), 则二次 $PE(\alpha)$ 方法收敛.

例 3.8 设 A 形如式 (3.137), 其中

$$B_i = \begin{bmatrix} 4 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 4 \end{bmatrix}, \quad A_i = C_i = \begin{bmatrix} -1 & & \\ & -1 & \\ & & -1 \end{bmatrix}.$$

可以验证 A 为对称正定矩阵. 取 $b = (1, 2, \dots, n)^T$, 用二次 $PE(\alpha)$ 方法求解线性方程组 $Ax = b$, 计算结果如表 3.8 所示 (终止准则为 10^{-12} , 时间单位为秒).

表 3.8 二次 $PE(\alpha)$ 方法的数值结果

参数 α	$n = 6000$			$n = 12000$		
	迭代次数	计算时间	相对残差	迭代次数	计算时间	相对残差
0.0	15	0.6491	2.6966e-13	15	1.2793	7.6464e-13
1.0	13	0.5943	1.5076e-13	13	1.1402	4.2746e-13
2.0	10	0.4680	5.2057e-13	11	0.9823	9.7413e-14
3.0	7	0.3516	1.3442e-13	7	0.6955	3.8024e-13
4.0	8	0.3927	7.4033e-13	9	0.8440	7.2741e-14
5.0	12	0.5265	8.9187e-14	12	1.0605	2.5292e-13
SBGS	29	1.1483	8.4202e-13	30	2.3955	9.4989e-13

表 3.8 中最后一行是用 SBGS 方法求解该方程组时的迭代次数、计算时间和相对残差. 易见, 当参数 $\alpha = 3.0$ 时, 二次 $PE(\alpha)$ 方法的收敛速度较快.

习题 3

3.1 若将方程组 $Ax = b$ 的每个方程两边除以 A 的相应对角元, 再对新得到的方程组构造 Richardson 迭代格式. 试证明它就是 $Ax = b$ 的 Jacobi 迭代格式.

3.2 设 $A \in \mathbf{R}^{n \times n}$ 对称正定, 其最小特征值和最大特征值分别为 λ_1 和 λ_n . 证明: 迭代法

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

收敛的充分必要条件是 $0 < \theta < 2/\lambda_n$.

3.3 设矩阵 $B \in \mathbb{R}^{n \times n}$ 满足 $\rho(B) = 0$. 证明: 对于任意的向量 $f, x^{(0)} \in \mathbb{R}^n$, 定常迭代格式

$$x^{(k+1)} = Bx^{(k)} + f, \quad k = 0, 1, \dots$$

最多迭代 n 次即可得到不动点方程 $x = Bx + f$ 的精确解.

3.4 设对称矩阵 $A = (a_{ij}) \in \mathbb{R}^{n \times n}$ 是非奇异的, 且 $a_{ii} > 0 (i = 1, 2, \dots, n)$. 证明: 若求解 $Ax = b$ 的 Gauss-Seidel 迭代法对任意初始向量 $x^{(0)}$ 都收敛, 则 A 必是对称正定矩阵.

3.5 证明: 若方程组 $Ax = b$ 的系数矩阵 A 是严格对角占优或不可约对角占优的, 且松弛参数 $0 < \omega < 2$, 则 SOR 迭代法是收敛的.

3.6 若存在对称正定矩阵 P , 使得

$$B = P - H^T P H$$

为对称正定矩阵. 证明: 迭代格式

$$x^{(k+1)} = Hx^{(k)} + f, \quad k = 0, 1, 2, \dots$$

收敛.

3.7 若存在可逆矩阵 P 使得 $P^{-1}(I - B)P$ 为对称正定矩阵, 则称单步定常迭代法

$$x^{(k+1)} = Bx^{(k)} + f, \quad k = 0, 1, \dots$$

是可对角化的. 试证明: 对于可对角化的单步定常迭代格式,

(1) 其迭代矩阵 B 的所有特征值都是小于 1 的实数;

(2) 存在实数 γ , 使得其外推格式

$$x^{(k+1)} = (1 - \gamma)x^{(k)} + \gamma(Bx^{(k)} + f), \quad k = 0, 1, \dots$$

收敛.

3.8 对于单步定常迭代格式

$$x^{(k+1)} = Bx^{(k)} + f, \quad k = 0, 1, \dots$$

若其迭代矩阵 $B \in \mathbb{R}^{n \times n}$ 的特征值 $\lambda_1, \lambda_2, \lambda_3$ 满足 $|\lambda_1| < 1, |\lambda_2| = 1, |\lambda_3| > 1$. 试问: 如何选取初始向量 $x^{(0)}$ 使得该迭代格式是收敛的或发散的?

3.9 证明: 严格对角占优矩阵不一定是不可约对角占优矩阵, 反之亦然. 试举例加以说明.

3.10 设 Hermite 矩阵 $A \in \mathbb{C}^{n \times n}$ 有分裂 $A = M - N$, 其中 M 非奇异. 试证明:

(1) $M^H + N$ 也是 Hermite 矩阵;

(2) 对于任意的 $x \in \mathbb{C}^n$, 都有

$$x^H Ax - \tilde{x}^H A \tilde{x} = (x - \tilde{x})^H (M^H + N)(x - \tilde{x}),$$

式中: $\tilde{x} = M^{-1}Nx$;

(3) 若 A 和 $M^H + N$ 均正定, 则 $\rho(M^{-1}N) < 1$;

(4) 分别对 Jacobi 迭代法和 SOR 迭代法, 写出 $M^H + N$ 的表达式.

第 4 章 线性方程组的 Krylov 子空间迭代法

本章着重介绍求解大型稀疏线性方程组的 Krylov 子空间方法. 考虑线性方程组

$$Ax = b, \quad (4.1)$$

式中: 矩阵 $A \in \mathbb{R}^{n \times n}$ 和向量 $b \in \mathbb{R}^n$ 是已经给定的, 而 $x \in \mathbb{R}^n$ 是待求的未知向量. 取初始向量 $x_0 \in \mathbb{R}^n$, 则解方程组 (4.1) 等价于求解

$$Az = r_0, \quad x = x_0 + z, \quad r_0 = b - Ax_0. \quad (4.2)$$

所谓求解线性方程组 (4.2) 的 Krylov 子空间方法, 就是先构造一个 Krylov 子空间

$$\mathcal{K}_k(A, r_0) = \text{span}\{r_0, Ar_0, \dots, A^{k-1}r_0\}; \quad (4.3)$$

然后设法求一个 $z_k \in \mathcal{K}_k(A, r_0)$, 使其在某种意义下是方程组 (4.2) 解的最佳逼近, 然后得到原方程组 (4.1) 的近似解 $x_k = x_0 + z_k$.

因此, 具体实现这一思想的关键是如何定量地刻画“最佳逼近”. 目前最常用的最佳性标准主要有两类: 一类是残量极小化标准; 另一类是残量正交化标准.

残量极小化标准是指求 $x_k \in x_0 + \mathcal{K}_k(A, r_0)$, 使得

$$\|r_k\|_2 = \min\{\|b - Ax\|_2 : x \in x_0 + \mathcal{K}_k(A, r_0)\}, \quad (4.4)$$

式中: $r_k = b - Ax_k$ 为 x_k 的残差向量 (简称残量).

本章主要介绍由这一标准导出的七种方法: 求解对称正定线性方程组的共轭梯度法 (CG 方法), 求解对称不定线性方程组的极小残量法 (MINRES 方法) 和 SYMMLQ 方法, 求解非对称线性方程组的广义极小残量法 (GMRES 方法)、拟极小残量法 (QMR 方法)、LSQR 方法和广义共轭残量法等.

残量正交化标准是指求 $x_k \in x_0 + \mathcal{K}_k(A, r_0)$, 使得

$$r_k = b - Ax_k \perp \mathcal{L}_k, \quad (4.5)$$

式中: \mathcal{L}_k 是另一选定的 k 维子空间, 称为约束空间.

\mathcal{L}_k 通常有三种典型取法: ① $\mathcal{L}_k = \mathcal{K}_k(A, r_0)$; ② $\mathcal{L}_k = A\mathcal{K}_k(A, r_0)$; ③ $\mathcal{L}_k = \mathcal{K}_k(A^T, r_0)$. 式 (4.5) 常称为 Galerkin 条件. 这里主要介绍这方面的三类方法: 双共轭梯度法 (BiCG 方法)、共轭梯度平方法 (CGS 方法) 和稳定化双共轭梯度法 (BICGSTAB 方法).

Krylov 子空间方法的一个最显著的特点就是整个计算过程只涉及矩阵 A 与某些向量的乘积, 这使得在算法执行过程中可以充分利用 A 的特殊性. 因此, 这种类型的线性方程组的求解方法就特别适用于大型稀疏线性方程组的求解. 当然, 对某些具有某种特殊结构的大型稠密线性方程组也是适用的.

4.1 共轭梯度法

本节介绍求解对称正定线性方程组的一类最著名的 Krylov 子空间方法——共轭梯度法.

4.1.1 基本 CG 方法

设 $A \in \mathbb{R}^{n \times n}$ 和 $b \in \mathbb{R}^n$ 已经给定, 其中 A 是对称正定的. 对于给定初始向量 $x_0 \in \mathbb{R}^n$, 求一个 $z \in \mathbb{R}^n$ 满足式 (4.2), 便可得到方程组 (4.1) 的解 $x = x_0 + z$.

由于 A 是对称正定的, 所以在 \mathbb{R}^n 上可以定义两种向量范数:

$$\|v\|_A = \sqrt{v^T A v} \quad \text{和} \quad \|v\|_{A^{-1}} = \sqrt{v^T A^{-1} v}, \quad v \in \mathbb{R}^n. \quad (4.6)$$

考虑由式 (4.2) 的系数矩阵 A 和右端项 r_0 所产生的 Krylov 子空间 $\mathcal{K}_k(A, r_0)$. 并记式 (4.1) 的真解为 x^* 即 $x^* = A^{-1}b$. 考虑求一个 $x_k \in x_0 + \mathcal{K}_k(A, r_0)$, 使其在某种条件下是 x^* 的最佳逼近. 下面的定理是关于这种最佳性的几种等价表述.

定理 4.1 符号如上所述, 则下面三条等价:

(1) 向量 $x_k \in x_0 + \mathcal{K}_k(A, r_0)$, 满足

$$\|x_k - x^*\|_A = \min\{\|x - x^*\|_A : x \in x_0 + \mathcal{K}_k(A, r_0)\} \quad (\text{误差极小}). \quad (4.7)$$

(2) 向量 $x_k \in x_0 + \mathcal{K}_k(A, r_0)$, 满足

$$\|b - Ax_k\|_{A^{-1}} = \min\{\|b - Ax\|_{A^{-1}} : x \in x_0 + \mathcal{K}_k(A, r_0)\} \quad (\text{残量极小}). \quad (4.8)$$

(3) 向量 $x_k \in x_0 + \mathcal{K}_k(A, r_0)$, 满足

$$z^T(b - Ax_k) = 0 \quad (\text{即 } r_k \perp \mathcal{K}_k(A, r_0)), \quad \forall z \in \mathcal{K}_k(A, r_0) \quad (\text{残量正交}). \quad (4.9)$$

注 4.1 定理 4.1 说明, 求 $x_k \in x_0 + \mathcal{K}_k(A, r_0)$ 使其在 $\|\cdot\|_A$ 范数下与真解 x^* 的距离最小就等价于其残量 $r_k = b - Ax_k$ 在 $\|\cdot\|_{A^{-1}}$ 范数下达到最小, 也等价于其残量 $r_k = b - Ax_k$ 与 $\mathcal{K}_k(A, r_0)$ 是正交的.

共轭梯度法基本迭代格式的推导需要利用对称 Lanczos 正交化过程. 假设已求得一个以 $v_1 = r_0/\|r_0\|_2$ 为初始向量的长度 k 的 Lanczos 分解 (见算法 2.13):

$$AV_k = V_k T_k + \beta_k v_{k+1} e_k^T, \quad (4.10)$$

式中: $V_k^T V_k = I_k$, $T_k = V_k^T A V_k$ 为如下形式的对称三对角矩阵

$$T_k = \begin{bmatrix} \alpha_1 & \beta_2 & & & \\ \beta_2 & \alpha_2 & \beta_3 & & \\ & \ddots & \ddots & \ddots & \\ & & \beta_{k-1} & \alpha_{k-1} & \beta_k \\ & & & \beta_k & \alpha_k \end{bmatrix},$$

且 $v_{k+1}^T V_k = 0$, $\|v_{k+1}\|_2 = 1$. 这里假定初始向量 $v_1 = r_0/\|r_0\|_2$, $r_0 = b - Ax_0$. 经过 Lanczos 正交化后, 得到的 $V_k = [v_1, v_2, \dots, v_k]$, 其列向量已构成了 Krylov 子空间 $\mathcal{K}_k(A, r_0)$ 的一组标准正交基, 即有

$$\mathcal{K}_k(A, r_0) = \text{span}\{r_0, Ar_0, \dots, A^{k-1}r_0\} = \mathcal{R}(V_k). \quad (4.11)$$

现在希望求一个 $\mathbf{x}_k \in \mathbf{x}_0 + \mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$ 使得式 (4.9) 成立. 由式 (4.11) 可知, 这就相当于求 $\mathbf{z}_k \in \mathbb{R}^k$, 使得

$$\mathbf{V}_k^T(\mathbf{r}_0 - \mathbf{A}\mathbf{V}_k\mathbf{z}_k) = \mathbf{V}_k^T(\mathbf{b} - \mathbf{A}\mathbf{x}_k) = \mathbf{0}, \quad (4.12)$$

式中: $\mathbf{x}_k = \mathbf{x}_0 + \mathbf{V}_k\mathbf{z}_k$.

注意到 $\mathbf{T}_k = \mathbf{V}_k^T \mathbf{A} \mathbf{V}_k$ 和 $\mathbf{r}_0 = \|\mathbf{r}_0\|_2 \mathbf{v}_1$, 由式 (4.12), 得

$$\mathbf{T}_k \mathbf{z}_k = \beta \mathbf{e}_1^{(k)}, \quad (4.13)$$

式中: $\beta = \|\mathbf{r}_0\|_2$, $\mathbf{e}_1^{(k)}$ 表示第 1 个分量为 1、其余分量为 0 的 k 维列向量.

注意到 \mathbf{A} 对称正定蕴涵着 \mathbf{T}_k 也是正定的, 便知 \mathbf{T}_k 的 \mathbf{LDL}^T 分解存在, 设其为

$$\mathbf{T}_k = \mathbf{L}_k \mathbf{D}_k \mathbf{L}_k^T, \quad (4.14)$$

式中: \mathbf{L}_k 为单位下三角矩阵; $\mathbf{D}_k = \text{diag}(\delta_1, \delta_2, \dots, \delta_k)$, $\delta_i > 0$, $i = 1, 2, \dots, k$. \mathbf{T}_k 的三对角结构蕴涵着 \mathbf{L}_k 必有如下形状:

$$\mathbf{L}_k = \begin{bmatrix} 1 & & & \\ \gamma_1 & 1 & & \\ & \ddots & \ddots & \\ & & \gamma_{k-1} & 1 \end{bmatrix}. \quad (4.15)$$

利用分解 (4.14), 线性方程组 (4.13) 的解可以表示为

$$\mathbf{z}_k = \mathbf{L}_k^{-T} \mathbf{D}_k^{-1} \mathbf{L}_k^{-1} (\beta \mathbf{e}_1^{(k)}),$$

从而有

$$\mathbf{x}_k = \mathbf{x}_0 + \mathbf{V}_k \mathbf{L}_k^{-T} \mathbf{D}_k^{-1} \mathbf{L}_k^{-1} (\beta \mathbf{e}_1^{(k)}). \quad (4.16)$$

再令

$$\tilde{\mathbf{P}}_k = [\tilde{\mathbf{p}}_1, \tilde{\mathbf{p}}_2, \dots, \tilde{\mathbf{p}}_k] = \mathbf{V}_k \mathbf{L}_k^{-T},$$

$$\mathbf{z}_k = (\zeta_1, \zeta_2, \dots, \zeta_k)^T = \mathbf{D}_k^{-1} \mathbf{L}_k^{-1} (\beta \mathbf{e}_1^{(k)}),$$

则式 (4.16) 又可写为

$$\mathbf{x}_k = \mathbf{x}_0 + \tilde{\mathbf{P}}_k \mathbf{z}_k. \quad (4.17)$$

注意到

$$\begin{aligned} \left[\begin{array}{c|c} \mathbf{T}_k & \beta \mathbf{e}_k \\ \hline \beta \mathbf{e}_k^T & \alpha_{k+1} \end{array} \right] &= \mathbf{T}_{k+1} := \mathbf{L}_{k+1} \mathbf{D}_{k+1} \mathbf{L}_{k+1}^T \\ &= \left[\begin{array}{c|c} \mathbf{L}_k & \mathbf{0} \\ \hline \gamma_k \mathbf{e}_k^T & 1 \end{array} \right] \left[\begin{array}{c|c} \mathbf{D}_k & \mathbf{0} \\ \hline \mathbf{0} & \delta_{k+1} \end{array} \right] \left[\begin{array}{c|c} \mathbf{L}_k & \mathbf{0} \\ \hline \gamma_k \mathbf{e}_k^T & 1 \end{array} \right]^T, \end{aligned}$$

式中: \mathbf{e}_k 表示第 k 个分量为 1、其余分量为 0 的 k 维列向量. 于是有

$$\begin{aligned}
 \tilde{\mathbf{P}}_{k+1} &= \mathbf{V}_{k+1} \mathbf{L}_{k+1}^{-\mathrm{T}} = [\mathbf{V}_k, \mathbf{v}_{k+1}] \begin{bmatrix} \mathbf{L}_k^{-\mathrm{T}} & \mathbf{0} \\ \gamma_k \mathbf{e}_k^{\mathrm{T}} & 1 \end{bmatrix}^{-\mathrm{T}} \\
 &= [\mathbf{V}_k, \mathbf{v}_{k+1}] \begin{bmatrix} \mathbf{L}_k^{-\mathrm{T}} & -\gamma_k \mathbf{L}_k^{-\mathrm{T}} \mathbf{e}_k^{(k)} \\ \mathbf{0} & 1 \end{bmatrix} \\
 &= [\mathbf{V}_k \mathbf{L}_k^{-\mathrm{T}}, \mathbf{v}_{k+1} - \gamma_k \mathbf{V}_k \mathbf{L}_k^{-\mathrm{T}} \mathbf{e}_k^{(k)}] \\
 &= [\tilde{\mathbf{P}}_k, \mathbf{v}_{k+1} - \gamma_k \tilde{\mathbf{P}}_k \mathbf{e}_k^{(k)}] := [\tilde{\mathbf{P}}_k, \tilde{\mathbf{p}}_{k+1}].
 \end{aligned}$$

由此可得

$$\tilde{\mathbf{p}}_{k+1} = \mathbf{v}_{k+1} - \gamma_k \tilde{\mathbf{P}}_k \mathbf{e}_k^{(k)} = \mathbf{v}_{k+1} - \gamma_k \tilde{\mathbf{p}}_k, \quad (4.18)$$

$$\begin{aligned}
 \mathbf{z}_{k+1} &= \mathbf{D}_{k+1}^{-1} \mathbf{L}_{k+1}^{-1} (\beta \mathbf{e}_1^{(k+1)}) \\
 &= \begin{bmatrix} \mathbf{D}_k^{-1} & \mathbf{0} \\ \mathbf{0} & \delta_{k+1}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{L}_k^{-1} & \mathbf{0} \\ -\gamma_k \mathbf{e}_k^{\mathrm{T}} \mathbf{L}_k^{-1} & 1 \end{bmatrix} (\beta \mathbf{e}_1^{(k+1)}) \\
 &= \begin{bmatrix} \mathbf{D}_k^{-1} \mathbf{L}_k^{-1} (\beta \mathbf{e}_1^{(k)}) \\ \zeta_{k+1} \end{bmatrix} = \begin{bmatrix} \mathbf{z}_k \\ \zeta_{k+1} \end{bmatrix},
 \end{aligned}$$

式中: $\zeta_{k+1} = -\beta \gamma_k \delta_{k+1}^{-1} \mathbf{e}_k^{\mathrm{T}} \mathbf{L}_k^{-1} \mathbf{e}_1^{(k)}$. 从而有

$$\mathbf{x}_{k+1} = \mathbf{x}_0 + \tilde{\mathbf{P}}_{k+1} \mathbf{z}_{k+1} = \mathbf{x}_0 + [\tilde{\mathbf{P}}_k, \tilde{\mathbf{p}}_{k+1}] \begin{bmatrix} \mathbf{z}_k \\ \zeta_{k+1} \end{bmatrix} = \mathbf{x}_k + \zeta_{k+1} \tilde{\mathbf{p}}_{k+1}. \quad (4.19)$$

由式 (4.19), 得

$$\begin{aligned}
 \mathbf{r}_{k+1} &= \mathbf{b} - \mathbf{A} \mathbf{x}_{k+1} = \mathbf{b} - \mathbf{A} \mathbf{x}_k - \zeta_{k+1} \mathbf{A} \tilde{\mathbf{p}}_{k+1} \\
 &= \mathbf{r}_k - \zeta_{k+1} \mathbf{A} \tilde{\mathbf{p}}_{k+1}.
 \end{aligned} \quad (4.20)$$

这样, 综合上面的推导, 得如下三个基本的迭代公式:

$$\begin{cases} \mathbf{x}_{k+1} = \mathbf{x}_k + \zeta_{k+1} \tilde{\mathbf{p}}_{k+1}, \\ \mathbf{r}_{k+1} = \mathbf{r}_k - \zeta_{k+1} \mathbf{A} \tilde{\mathbf{p}}_{k+1}, \\ \tilde{\mathbf{p}}_{k+1} = \mathbf{v}_{k+1} - \gamma_k \tilde{\mathbf{p}}_k. \end{cases} \quad (4.21)$$

下面设法消去式 (4.21) 第 3 式中的 \mathbf{v}_{k+1} . 由 $\mathbf{r}_0 \in \mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$, $\mathbf{x}_k \in \mathbf{x}_0 + \mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$ 可知, $\mathbf{r}_k = \mathbf{b} - \mathbf{A} \mathbf{x}_k \in \mathcal{K}_{k+1}(\mathbf{A}, \mathbf{r}_0)$, 从而 \mathbf{r}_k 可以由 \mathbf{V}_{k+1} 的列向量线性表出, 即

$$\mathbf{r}_k = \tilde{\gamma}_1 \mathbf{v}_1 + \tilde{\gamma}_2 \mathbf{v}_2 + \cdots + \tilde{\gamma}_k \mathbf{v}_k + \tilde{\gamma}_{k+1} \mathbf{v}_{k+1}.$$

又由于 $\mathbf{V}_k^T \mathbf{r}_k = \mathbf{0}$, 故 $\tilde{\gamma}_i = 0, i = 1, 2, \dots, k$, 从而

$$\mathbf{r}_k = \tilde{\gamma}_{k+1} \mathbf{v}_{k+1}. \quad (4.22)$$

由式 (4.22) 可知 $|\tilde{\gamma}_{k+1}| = \|\mathbf{r}_k\|_2$, 不妨设 $\tilde{\gamma}_{k+1} = \|\mathbf{r}_k\|_2$, 并定义

$$\mathbf{p}_{k+1} = \|\mathbf{r}_k\|_2 \tilde{\mathbf{p}}_{k+1}, \quad (4.23)$$

则由式 (4.21) 和式 (4.22), 得

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k + \frac{\zeta_{k+1}}{\|\mathbf{r}_k\|_2} \mathbf{p}_{k+1}, \\ \mathbf{r}_{k+1} &= \mathbf{r}_k - \frac{\zeta_{k+1}}{\|\mathbf{r}_k\|_2} \mathbf{A} \mathbf{p}_{k+1}, \\ \mathbf{p}_{k+1} &= \mathbf{r}_k - \frac{\|\mathbf{r}_k\|_2 \gamma_k}{\|\mathbf{r}_{k-1}\|_2} \mathbf{p}_k. \end{aligned}$$

令

$$\alpha_{k+1} = \frac{\zeta_{k+1}}{\|\mathbf{r}_k\|_2}, \quad \beta_k = -\frac{\|\mathbf{r}_k\|_2 \gamma_k}{\|\mathbf{r}_{k-1}\|_2},$$

即得

$$\begin{cases} \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_{k+1} \mathbf{p}_{k+1}, \\ \mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_{k+1} \mathbf{A} \mathbf{p}_{k+1}, \\ \mathbf{p}_{k+1} = \mathbf{r}_k + \beta_k \mathbf{p}_k. \end{cases} \quad (4.24)$$

下面导出不涉及分解式 (4.14) 的信息计算 α_{k+1} 和 β_k 的公式. 为此, 先导出向量组 $\{\mathbf{p}_i\}$ 和 $\{\mathbf{r}_i\}$ 所具有的基本性质.

由 $\tilde{\mathbf{P}}_k$ 的定义, 可以立即导出

$$\tilde{\mathbf{P}}_k^T \mathbf{A} \tilde{\mathbf{P}}_k = \mathbf{L}_k^{-1} \mathbf{V}_k^T \mathbf{A} \mathbf{V}_k \mathbf{L}_k^{-T} = \mathbf{D}_k.$$

这表明 $\tilde{\mathbf{p}}_i$ 之间是 \mathbf{A} 正交的 (或称为 \mathbf{A} 共轭的), 即有

$$\tilde{\mathbf{p}}_i^T \mathbf{A} \tilde{\mathbf{p}}_j = 0, \quad i \neq j.$$

于是有

$$\mathbf{p}_i^T \mathbf{A} \mathbf{p}_j = 0, \quad i \neq j. \quad (4.25)$$

此外, 式 (4.1.1) 的第 1 式和式 (4.23) 蕴涵着

$$\mathcal{K}_k(\mathbf{A}, \mathbf{r}_0) = \mathcal{R}(\mathbf{V}_k) = \mathcal{R}(\tilde{\mathbf{P}}_k) = \text{span}\{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_k\}, \quad (4.26)$$

即 $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_k$ 是 $\mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$ 的一组 \mathbf{A} 正交基, 再由式 (4.22) 和 \mathbf{v}_i 的相互正交性, 又可以导出

$$\mathbf{r}_i^T \mathbf{r}_j = 0, \quad i \neq j. \quad (4.27)$$

即残量是相互正交的. 这样, 再由 $\mathcal{K}_k(\mathbf{A}, \mathbf{r}_0) = \mathcal{R}(\mathbf{V}_k)$ 便有

$$\mathcal{K}_k(\mathbf{A}, \mathbf{r}_0) = \text{span}\{\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_{k-1}\}, \quad (4.28)$$

即 $\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_{k-1}$ 构成了 $\mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$ 的一组正交基.

在式 (4.24) 的第 3 式两边左乘 $\mathbf{p}_k^T \mathbf{A}$ 并利用式 (4.25), 得

$$0 = \mathbf{p}_k^T \mathbf{A} \mathbf{p}_{k+1} = \mathbf{p}_k^T \mathbf{A} \mathbf{r}_k + \beta_k \mathbf{p}_k^T \mathbf{A} \mathbf{p}_k.$$

于是有

$$\beta_k = -\frac{\mathbf{p}_k^T \mathbf{A} \mathbf{r}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k}. \quad (4.29)$$

再在式 (4.24) 的第 2 式两边左乘 \mathbf{r}_k^T 并利用式 (4.27), 得

$$0 = \mathbf{r}_k^T \mathbf{r}_{k+1} = \mathbf{r}_k^T \mathbf{r}_k - \alpha_{k+1} \mathbf{r}_k^T \mathbf{A} \mathbf{p}_{k+1}.$$

因此又有

$$\alpha_{k+1} = \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_k^T \mathbf{A} \mathbf{p}_{k+1}}. \quad (4.30)$$

这样就得到了不需要分解式 (4.14) 的信息就可计算式 (4.24) 中的系数 β_k 和 α_{k+1} 的公式.

事实上, 利用式 (4.25) 和式 (4.27) 还可导出更有效的计算公式. 首先在式 (4.24) 的第 3 式两边左乘 $\mathbf{p}_{k+1}^T \mathbf{A}$, 并利用式 (4.25), 得

$$\mathbf{p}_{k+1}^T \mathbf{A} \mathbf{p}_{k+1} = \mathbf{p}_{k+1}^T \mathbf{A} \mathbf{r}_k = \mathbf{r}_k^T \mathbf{A} \mathbf{p}_{k+1}. \quad (4.31)$$

将式 (4.24) 第 2 式中的 k 用 $k-1$ 替换, 有

$$\mathbf{r}_k = \mathbf{r}_{k-1} - \alpha_k \mathbf{A} \mathbf{p}_k.$$

将上式两边分别左乘 \mathbf{r}_k^T 和 \mathbf{r}_{k-1}^T , 并利用式 (4.27), 得

$$\mathbf{r}_k^T \mathbf{r}_k = -\alpha_k \mathbf{r}_k^T \mathbf{A} \mathbf{p}_k, \quad (4.32)$$

$$\mathbf{r}_{k-1}^T \mathbf{r}_{k-1} = \alpha_k \mathbf{r}_{k-1}^T \mathbf{A} \mathbf{p}_k = \alpha_k \mathbf{p}_k^T \mathbf{A} \mathbf{p}_k, \quad (4.33)$$

这里式 (4.33) 的第 2 个等号利用了式 (4.31).

现将式 (4.31) 代入式 (4.30), 得

$$\alpha_{k+1} = \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{p}_{k+1}^T \mathbf{A} \mathbf{p}_{k+1}}. \quad (4.34)$$

再将式 (4.32)、式 (4.33) 与式 (4.29) 相结合, 有

$$\beta_k = \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_{k-1}^T \mathbf{r}_{k-1}}. \quad (4.35)$$

显然, 这两个计算公式比原来的计算公式减少了内积和矩阵乘向量的次数, 因此更有效.

综合上述讨论, 可得共轭梯度法 (简称 CG 方法) 求解对称正定线性方程组的基本迭代格式如下.

选取初值 x_0 , 计算 $r_0 = b - Ax_0$; $p_1 = r_0$; $k = 0$;

while ($\|r_k\|_2 / \|r_0\|_2 > \varepsilon$)

$k = k + 1$;

$\alpha_k = \frac{(r_{k-1}, r_{k-1})}{(p_k, Ap_k)}$; $x_k = x_{k-1} + \alpha_k p_k$;

$r_k = r_{k-1} - \alpha_k Ap_k$;

$\beta_k = \frac{(r_k, r_k)}{(r_{k-1}, r_{k-1})}$; $p_{k+1} = r_k + \beta_k p_k$;

end

为了优化算法程序 (考虑计算量和存储量), 通常写成如下形式.

算法 4.1 (CG 方法) 给定对称正定方程组 $Ax = b$ 和 $\varepsilon > 0$. 本算法计算 x_k , 使得 $\|r_k\|_2 / \|r_0\|_2 \leq \varepsilon$, 其中 $r_k = b - Ax_k$.

选取 x_0 ; $r_0 = b - Ax_0$; $p_1 = r_0$; $\rho_0 = r_0^T r_0$; $k = 0$;

while ($\|r_k\|_2 / \|r_0\|_2 > \varepsilon$)

$k = k + 1$;

$z_k = Ap_k$; $\alpha_k = \rho_{k-1} / z_k^T p_k$;

$x_k = x_{k-1} + \alpha_k p_k$; $r_k = r_{k-1} - \alpha_k z_k$;

$\rho_k = r_k^T r_k$; $\beta_k = \rho_k / \rho_{k-1}$;

$p_{k+1} = r_k + \beta_k p_k$;

end

算法 4.1 的迭代终止准则也可以用 $\rho_k \leq \varepsilon \rho_0$, 其中 $\varepsilon > 0$ 是给定的允许误差. CG 方法的 MATLAB 程序如下:

%共轭梯度法程序-mcg.m

function [x,iter,time,res,resvec]=mcg(A,b,x,max_it,tol)

%输入:系数阵A,右端量b,初始值x,容许误差tol,最大迭代数max_it

%输出:解向量x,迭代次数iter,CPU时间time,

% 终止时相对残差模res,相对残差模向量resvec

tic;r=b-A*x;p=r;rho=r'*r;

mr=sqrt(rho);iter=0;

while (iter<max_it)

iter=iter+1;

z=A*p; alpha=rho/(z'*p);

```

x=x+alpha*p; r=r-alpha*z;
rho1=r'*r; beta=rho1/rho;
p=r+beta*p;
res=sqrt(rho1)/mr;
resvec(iter)=res;
if (res<tol), break; end
rho=rho1;
end
time=toc;

```

例 4.1 假定线性方程组 $Ax = b$ 的系数矩阵 A 和右端项 b 分别为

$$A = \begin{bmatrix} 1 & -1 & & & \\ -1 & 2 & \ddots & & \\ & \ddots & \ddots & \ddots & \\ & & & -1 & n \end{bmatrix}, \quad b = A \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}.$$

显然, 此方程组的真解为 $x^* = (1, 1, \dots, 1)^T$. 应用共轭梯度法求解该线性方程组, 取 $n = 1000$, 迭代 193 步后得到的近似解 $\hat{x} = x_{193}$ 满足

$$\|\hat{x} - x^*\|_2 = 3.7417 \times 10^{-8},$$

迭代过程的收敛轨迹如图 4.1 所示, 其中横坐标为迭代步数 k , 纵坐标为相对残差 $\|r_k\|_2 / \|r_0\|_2$, 这里 r_k 是第 k 步得到的残差.

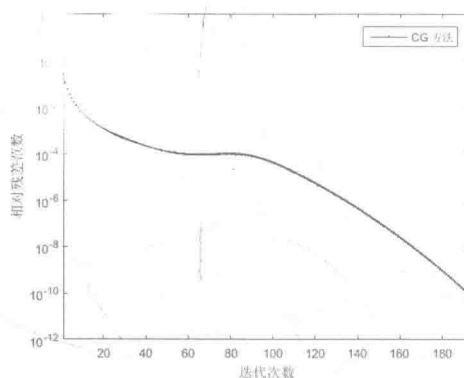


图 4.1 CG 算法的收敛特性

4.1.2 收敛性分析

在没有误差的情况下, 由式 (4.27) 可知, 必存在一个 $\ell \leq n$, 使得 $r_\ell = 0$, 即有 $x^* = x_\ell$. 换言之, 算法在有限步可得到方程组 (4.1) 的精确解. 从这个意义上而言, 共

轭梯度法是一种直接法. 但实际使用时, 一般是将其作为一种迭代法来用的, 其原因有二: 一是在有误差的计算机上, 有限步终止的结论已经不再成立; 二是即使出现了 $r_\ell = 0$, 但如果 ℓ 非常大, 也是难以忍受的, 因为实际应用中所考虑的问题其阶数 n 是十分巨大的.

下面介绍共轭梯度法 (作为一种迭代法) 的收敛性和收敛速度. 由于收敛性分析需要利用 Chebyshev 多项式的性质. 在第 3 章已给出 k 次 Chebyshev 多项式 $C_k(t)$ 的定义:

$$C_k(t) = \begin{cases} \cos(k \arccos t), & |t| \leq 1, \\ \cosh(k \operatorname{arccosh} t), & |t| > 1. \end{cases} \quad (4.36)$$

下面再回顾一下 Chebyshev 多项式的一些最常用的性质.

定理 4.2 设 $C_k(t)$ 是由式 (4.36) 定义的 Chebyshev 多项式, 则它具有如下性质:

(1) $C_k(t)$ 具有递推式

$$C_0(t) = 1; C_1(t) = t; C_{k+1}(t) = 2tC_k(t) - C_{k-1}(t), \quad k = 1, 2, \dots$$

(2) 对于 $|t| > 1$, $C_k(t)$ 有表达式

$$C_k(t) = \frac{1}{2} \left[\left(t + \sqrt{t^2 - 1} \right)^k + \left(t + \sqrt{t^2 - 1} \right)^{-k} \right]. \quad (4.37)$$

(3) 对任意的 $t \in [-1, 1]$, 有 $|C_k(t)| \leq 1$, 且

$$C_k(t_i^{(k)}) = (-1)^i, \quad t_i^{(k)} = \cos \frac{i\pi}{k}, \quad i = 0, 1, \dots, k,$$

即 $C_k(t)$ 在 $[-1, 1]$ 上刚好有 $k+1$ 个极值点, 在这些点上交错地取 1 和 -1.

(4) 对任意的 $s > 1$, 有

$$\min_{p \in \mathcal{P}_k, p(s)=1} \max_{t \in [-1, 1]} |p(t)| = \frac{1}{C_k(s)}, \quad (4.38)$$

式中:

$$\mathcal{P}_k = \{p: p \text{ 是次数不超过 } k \text{ 的实系数多项式}\},$$

而且上述极小极大问题有唯一解

$$p(t) = \frac{C_k(t)}{C_k(s)}. \quad (4.39)$$

证明 (1) 对于 $|t| \leq 1$, 令 $\theta = \arccos t$, 利用余弦函数的和差化积公式

$$\cos \alpha + \cos \beta = 2 \cos \frac{\alpha + \beta}{2} \cos \frac{\alpha - \beta}{2},$$

有

$$C_{k+1}(t) + C_{k-1}(t) = \cos(k+1)\theta + \cos(k-1)\theta = 2 \cos k\theta \cos \theta = 2tC_k(t).$$

对于 $|t| > 1$, 令 $\theta = \operatorname{arccosh}(t)$, 利用双曲余弦函数的定义

$$\cosh(\theta) = \frac{e^\theta + e^{-\theta}}{2},$$

有

$$\begin{aligned} C_{k+1}(t) + C_{k-1}(t) &= \cosh(k+1)\theta + \cosh(k-1)\theta \\ &= \frac{e^{(k+1)\theta} + e^{-(k+1)\theta}}{2} + \frac{e^{(k-1)\theta} + e^{-(k-1)\theta}}{2} \\ &= e^{k\theta} \cdot \frac{e^\theta + e^{-\theta}}{2} + e^{-k\theta} \cdot \frac{e^\theta + e^{-\theta}}{2} \\ &= 2 \cdot \frac{e^\theta + e^{-\theta}}{2} \cdot \frac{e^{k\theta} + e^{-k\theta}}{2} \\ &= 2 \cosh \theta \cosh(k\theta) = 2t C_k(t). \end{aligned}$$

(2) 对于 $|t| > 1$, 令 $\theta = \operatorname{arccosh}(t)$, 有 $t = \cosh(\theta) = e^\theta + e^{-\theta}/2$, 即

$$(e^\theta)^2 - 2te^\theta + 1 = 0.$$

得

$$e^\theta = \frac{2t + \sqrt{4t^2 - 4}}{2} = t + \sqrt{t^2 - 1}.$$

于是, 有

$$\begin{aligned} C_k(t) &= \cosh(k\theta) = \frac{1}{2}(e^{k\theta} + e^{-k\theta}) = \frac{1}{2}[(e^\theta)^k + (e^\theta)^{-k}] \\ &= \frac{1}{2} \left[\left(t + \sqrt{t^2 - 1} \right)^k + \left(t + \sqrt{t^2 - 1} \right)^{-k} \right]. \end{aligned}$$

(3) $\forall t \in [-1, 1]$, $|C_k(t)| \leq 1$ 成立是显然的. 直接验证, 得

$$C_k(t_i^{(k)}) = \cos(k \arccos t_i^{(k)}) = \cos(i\pi) = (-1)^i.$$

(4) 注意到 $s > 1$, $p(t)$ 是 k 次多项式且 $p(s) = 1$. 记

$$\alpha = 1/C_k(s) > 0, \quad p_{\max} = \max_{-1 \leq t \leq 1} |p(t)|.$$

若 $p_{\max} < \alpha$, 则

$$\begin{aligned} (-1)^i (\alpha C_k(t_i^{(k)}) - p(t_i^{(k)})) &= (-1)^i \alpha (-1)^i - (-1)^i p(t_i^{(k)}) \\ &= \alpha - (-1)^i p(t_i^{(k)}) \geq \alpha - p_{\max} > 0. \end{aligned}$$

这表明多项式 $\psi = \alpha C_k(t) - p(t)$ 在 $t_i^{(k)}$ ($i = 0, 1, \dots, k$) 处有交替的正负号. 所以它在 $[-1, 1]$ 中有 k 个零点. 但 $s \notin [-1, 1]$ 也是 k 次多项式 $\alpha C_k(t) - p(t)$ 的零点, 从而 k 次多项式 $\alpha C_k(t) - p(t)$ 有 $k+1$ 个零点, 矛盾. 故必有

$$\max_{-1 \leq t \leq 1} |p(t)| \geq \alpha = \frac{1}{C_k(s)}.$$

若取 $p(t) = \frac{C_k(t)}{C_k(s)}$, 则

$$\max_{-1 \leq t \leq 1} |p(t)| = \max_{-1 \leq t \leq 1} \left| \frac{C_k(t)}{C_k(s)} \right| = \frac{\max_{-1 \leq t \leq 1} |C_k(t)|}{C_k(s)} = \frac{1}{C_k(s)}.$$

由此可得

$$\min_{p \in \mathcal{P}_k, p(1)=1} \max_{t \in [-1, 1]} |p(t)| = \frac{1}{C_k(s)}.$$

证毕. □

推论 4.1 设 $a < b < 1$ (或 $1 < a < b$). 若 $p \in \mathcal{P}_k$ 且 $p(1) = 1$. 则

$$\min_{p \in \mathcal{P}_k, p(1)=1} \max_{t \in [a, b]} |p(t)| = \frac{1}{C_k(w(1))}, \tag{4.40}$$

式中:

$$w(t) = \begin{cases} \frac{2t - a - b}{b - a}, & a < b < 1, \\ \frac{2t - a - b}{a - b}, & 1 < a < b. \end{cases}$$

且上述极小极大问题的唯一解为

$$p(t) = \frac{C_k(w(t))}{C_k(w(1))}.$$

现在考虑 CG 方法的收敛性分析. 由定理 4.1 和前面的推导过程可知, 共轭梯度法第 k 步得到的近似解 \mathbf{x}_k 满足

$$\|\mathbf{x}_k - \mathbf{x}^*\|_A = \min \{ \|\mathbf{x} - \mathbf{x}^*\|_A : \mathbf{x} \in \mathbf{x}_0 + \mathcal{K}_k(\mathbf{A}, \mathbf{r}_0) \}.$$

由 Krylov 子空间的性质可知, $\mathbf{x} \in \mathbf{x}_0 + \mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$ 的充分必要条件是存在 $\varphi \in \mathcal{P}_{k-1}$ 使得 $\mathbf{x} = \mathbf{x}_0 + \varphi(\mathbf{A})\mathbf{r}_0$, 这里 \mathcal{P}_{k-1} 表示次数不超过 $k-1$ 的多项式的全体. 于是, 对任意的 $\mathbf{x} \in \mathbf{x}_0 + \mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$, 有

$$\mathbf{x}^* - \mathbf{x} = \mathbf{A}^{-1}\mathbf{b} - \mathbf{x}_0 - \varphi(\mathbf{A})\mathbf{r}_0 = \mathbf{A}^{-1}\psi(\mathbf{A})\mathbf{r}_0,$$

式中: $\psi(t) = 1 - t\varphi(t)$.

设 \mathbf{A} 的谱分解为 $\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$, 其中 $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n]$ 是正交矩阵, $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$, $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n > 0$. 现将 \mathbf{r}_0 按 \mathbf{A} 的特征向量展开, 有

$$\mathbf{r}_0 = \eta_1 \mathbf{u}_1 + \eta_2 \mathbf{u}_2 + \dots + \eta_n \mathbf{u}_n = \sum_{i=1}^n \eta_i \mathbf{u}_i.$$

这样, 便有

$$\|\mathbf{x}^* - \mathbf{x}\|_A^2 = \|\mathbf{A}^{-1}\psi(\mathbf{A})\mathbf{r}_0\|_A^2 = (\mathbf{A}^{-1}\psi(\mathbf{A})\mathbf{r}_0, \psi(\mathbf{A})\mathbf{r}_0)$$

$$\begin{aligned}
&= \left(\sum_{i=1}^n \eta_i A^{-1} \psi(A) u_i, \sum_{j=1}^n \eta_j \psi(A) u_j \right) \\
&= \left(\sum_{i=1}^n \eta_i \psi(\lambda_i) \lambda_i^{-1} u_i, \sum_{j=1}^n \eta_j \psi(\lambda_j) u_j \right) \\
&= \sum_{i=1}^n \eta_i^2 \psi^2(\lambda_i) \lambda_i^{-1} \leq \max_{1 \leq i \leq n} \psi^2(\lambda_i) \sum_{i=1}^n \lambda_i^{-1} \eta_i^2 \\
&= \max_{1 \leq i \leq n} \psi^2(\lambda_i) \|x^* - x_0\|_A^2,
\end{aligned}$$

从而有

$$\frac{\|x^* - x\|_A}{\|x^* - x_0\|_A} \leq \min_{\psi \in \mathcal{P}_k^0} \max_{1 \leq i \leq n} |\psi(\lambda_i)|, \quad (4.41)$$

式中: $\mathcal{P}_k^0 = \{\psi \in \mathcal{P}_k : \psi(0) = 1\}$. 令

$$a = \min_{1 \leq i \leq n} \lambda_i = \lambda_n, \quad b = \max_{1 \leq i \leq n} \lambda_i = \lambda_1,$$

则由推论 4.1, 知极小极大问题

$$\min_{\psi \in \mathcal{P}_k^0} \max_{\lambda \in [a, b]} |\psi(\lambda)|$$

有唯一的解

$$\hat{\psi}(\lambda) = \frac{C_k\left(\frac{b+a-2\lambda}{b-a}\right)}{C_k\left(\frac{b+a}{b-a}\right)},$$

而且有

$$\min_{\psi \in \mathcal{P}_k^0} \max_{\lambda \in [a, b]} |\psi(\lambda)| = \frac{1}{C_k\left(\frac{b+a}{b-a}\right)}, \quad (4.42)$$

这里假设了 $1 < a < b$. 将式 (4.41) 和式 (4.42) 相结合, 可证明如下结果.

定理 4.3 共轭梯度法产生的近似解 x_k 满足

$$\frac{\|x_k - x^*\|_A}{\|x_0 - x^*\|_A} \leq \frac{1}{C_k\left(1 + \frac{2}{\kappa - 1}\right)}, \quad (4.43)$$

式中: C_k 为 k 阶 Chebyshev 多项式, $\kappa = \kappa_2(A) = \|A\|_2 \cdot \|A^{-1}\|_2 = b/a$.

注 4.2 由定理 4.2 的结论 (2), 得

$$\left[C_k \left(1 + \frac{2}{\kappa - 1} \right) \right]^{-1} = \left[C_k \left(\frac{b+a}{b-a} \right) \right]^{-1} = \frac{2\sigma^k}{1 + \sigma^{2k}}, \quad (4.44)$$

式中: $\sigma = (\sqrt{b} - \sqrt{a})/(\sqrt{b} + \sqrt{a})$. 由式 (4.44) 可证

$$\left[C_k \left(1 + \frac{2}{\kappa - 1} \right) \right]^{-1} \leq 2 \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k. \quad (4.45)$$

由此可知, 系数矩阵的条件数 κ 越小, 共轭梯度法收敛就越快.

4.1.3 预处理 CG 方法

由式 (4.43), 在实际使用共轭梯度法时, 通常需对方程组 (4.1) 作预处理, 使其系数矩阵的谱相对集中, 这就是所谓的预处理共轭梯度法 (PCG 方法). 简单来说, 预处理共轭梯度法就是预先选择一个适当的对称正定矩阵 M , 使得矩阵 $M^{-1}A$ 的谱相对集中. 设 M 有分解 $M = LL^T$, 则方程组 (4.1) 经过预处理后变为

$$(L^{-1}AL^{-T})\tilde{x} = L^{-1}b, \quad x = L^{-T}\tilde{x}. \quad (4.46)$$

该方程组的系数矩阵 $L^{-1}AL^{-T}$ 仍为对称正定矩阵. 通过 M 和 L 的选取, 可以使得 $\tilde{A} = L^{-1}AL^{-T}$ 比 A 有更好的条件数. 所以预处理后的方程组 (4.46) 比原方程组更容易计算. 现在对方程组 (4.46) 使用共轭梯度法:

取 \tilde{x}_0 , 令 $\tilde{r}_0 = L^{-1}b - (L^{-1}AL^{-T})\tilde{x}_0$, $\tilde{p}_1 = \tilde{r}_0$, 对 $k = 1, 2, \dots$, 计算出

$$\begin{aligned} \tilde{\alpha}_k &= \frac{\tilde{r}_{k-1}^T \tilde{r}_{k-1}}{\tilde{p}_k^T \tilde{A} \tilde{p}_k}, \quad \tilde{x}_k = \tilde{x}_{k-1} + \tilde{\alpha}_k \tilde{p}_k, \quad \tilde{r}_k = \tilde{r}_{k-1} - \tilde{\alpha}_k \tilde{A} \tilde{p}_k, \\ \tilde{\beta}_k &= \frac{\tilde{r}_k^T \tilde{r}_k}{\tilde{r}_{k-1}^T \tilde{r}_{k-1}}, \quad \tilde{p}_{k+1} = \tilde{r}_k + \tilde{\beta}_k \tilde{p}_k. \end{aligned}$$

为了在公式中使用原来变量, 定义 $x_k = L^{-T}\tilde{x}_k$, 则 $r_k = b - Ax_k = L\tilde{r}_k$. 再令 $p_k = L^{-T}\tilde{p}_k$, $z_k = M^{-1}r_k$, 就得到了预处理的共轭梯度法, 具体步骤如下:

步 1, 给定 $x_0 \in \mathbb{R}^n$. 计算 $r_0 = b - Ax_0$, $z_0 = M^{-1}r_0$, $p_1 = z_0$.

步 2, 对 $k = 1, 2, \dots$, 计算

$$\begin{aligned} \alpha_k &= \frac{z_{k-1}^T r_{k-1}}{p_k^T A p_k}, \quad x_k = x_{k-1} + \alpha_k p_k, \\ r_k &= r_{k-1} - \alpha_k A p_k, \quad z_k = M^{-1}r_k, \end{aligned} \quad (4.47)$$

$$\beta_k = \frac{z_k^T r_k}{z_{k-1}^T r_{k-1}}, \quad p_{k+1} = z_k + \beta_k p_k. \quad (4.48)$$

为了优化算法程序, 写成如下形式.

算法 4.2 (PCG 方法) 给定对称正定方程组 $Ax = b$, 预处理子 M 和 $\varepsilon > 0$. 本算法计算 x_k , 使得 $\|r_k\|_2 / \|r_0\|_2 \leq \varepsilon$, 其中 $r_k = b - Ax_k$.

选取 x_0 ; $r_0 = b - Ax_0$; $z_0 = M^{-1}r_0$;

$p_1 = z_0$; $\rho_0 = z_0^T r_0$; $k = 0$;

while ($\|r_k\|_2 / \|r_0\|_2 > \varepsilon$)

$k = k + 1$;

$u_k = A p_k$; $\alpha_k = \rho_{k-1} / u_k^T p_k$;

$x_k = x_{k-1} + \alpha_k p_k$; $r_k = r_{k-1} - \alpha_k u_k$;

$$z_k = M^{-1}r_k; \rho_k = z_k^T r_k; \beta_k = \rho_k / \rho_{k-1};$$

$$p_{k+1} = z_k + \beta_k p_k;$$

end

算法 4.2 仅与预处理矩阵 M 有关, 而与 L 无关. 当 M 为单位阵时, 它就是没有经过预处理的共轭梯度法. 与共轭梯度法相比较, 预处理共轭梯度法仅增加了式 (4.47) 的工作量, 即每一步需要求一个以预处理矩阵 M 为系数矩阵的方程组得到 z_k . 算法 4.2 每一步的主要工作量是作一次矩阵 A 与向量的乘法操作和一次 M^{-1} 与向量的乘法操作.

下面讨论预处理矩阵 $M = LL^T$ 的选取. M 取法之一是 A 的不完全 Cholesky 分解, 即 $A = \hat{L}\hat{L}^T + R$, 然后取 $M := \hat{L}\hat{L}^T$, 其中 $R = A - M$ 满足某种稀疏模式, 使得它比较小.

如果对 A 进行分裂, 即 $A = M - N$, 则需要判断古典迭代格式中的 M 是否满足要求, 即是否是对称正定阵. 对于 Jacobi 迭代, M 为 A 的对角元所组成的对角矩阵, 它显然是对称正定的, 因此可以作为预处理矩阵. 这样的预处理称为 Jacobi 预处理. 对于 Gauss-Seidel 和 SOR 迭代, 其分裂格式中的 M 都是下三角矩阵, 不是对称正定阵, 因此不能作为预处理矩阵. 而对于 SSOR 迭代, 其分裂格式中的 M 是对称正定的, 因此可以作为预处理矩阵, 称其为 SSOR 预处理矩阵.

给出预处理共轭梯度法的 MATLAB 程序如下:

```
%预处理共轭梯度法程序-pcg.m
function [x,iter,time,res,resvec]=pcg(A,b,x,M,max_it,tol)
%输入:系数矩阵A,右端向量b,初始向量x,预处理子M,
%    容许误差tol,最大迭代次数max_it
%输出:解向量x,迭代次数iter,CPU时间time,
%    终止时相对残差模res,相对残差模向量resvec
tic; r=b-A*x; z=M\r; p=z;
rho=z'*r; mr=norm(r); iter=0;
while (iter<max_it)
    iter=iter+1;
    u=A*p; alpha=rho/(p'*u);
    x=x+alpha*p; r=r-alpha*u;
    z=M\r; rho1=z'*r;
    beta=rho1/rho; p=z+beta*p;
    res=norm(r)/mr;
    resvec(iter)=res;
    if (res<tol), break; end
    rho=rho1;
end
time=toc;
```

为了对预处理的效果有一个更直观更深刻的印象,下面再给出一个数值算例.

例 4.2 考虑例 4.1 中的线性方程组, 如果取预处理矩阵 M 为 A 的对角元构成的对角矩阵, 即 $M = \text{diag}(1, 2, \dots, n)$. 实验中取 $n = 1000$, 则预处理共轭梯度法在 12 步后得到的近似解 $\hat{x} = x_{12}$ 就满足 $\|\hat{x} - x^*\|_2 = 3.7305 \times 10^{-9}$. 迭代的收敛轨迹如图 4.2 所示, 其中横坐标表示迭代步数 k , 纵坐标表示相对残差 $\|r_k\|_2 / \|r_0\|_2$, 这里 r_k 是第 k 步得到的残差.

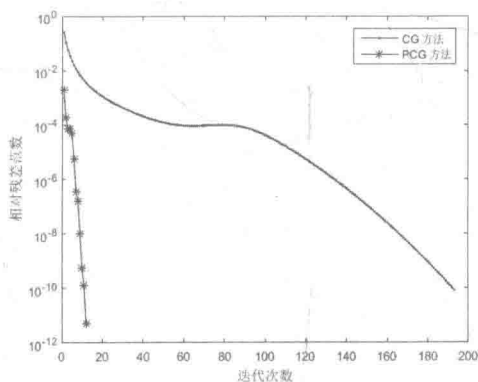


图 4.2 PCG 法和 CG 法的收敛特性

与例 4.1 的数值结果相比较可知, 虽然这里的预处理矩阵 M 选得非常简单, 但其加速效果是显著的. 为什么会这样呢? 因为预处理后矩阵 $M^{-1}A$ 的特征值分布很集中, 故其数值效果非常好.

4.1.4 CGNR 方法和 CGNE 方法

1. CGNR 方法

任何对称不定或非对称非奇异的方程组 $Ax = b$ 都可转化为一个对称正定的方程组 (即法方程)

$$A^T A x = A^T b. \quad (4.49)$$

对方程组 (4.49) 应用 CG 法, 就导出了下面的关于法方程残差的共轭梯度法 (CGNR).

算法 4.3 (CGNR 方法) 给定系数矩阵 A , 右端向量 b , 初始向量 x_0 和容许误差限 $\varepsilon > 0$. 本算法计算 x_k , 使得 $\|r_k\|_2 / \|r_0\|_2 \leq \varepsilon$, 其中 $r_k = b - Ax_k$.

选取 x_0 ; 计算 $r_0 = b - Ax_0$; $p_0 = A^T r_0$; $k = 0$;

while ($\|r_k\|_2 / \|r_0\|_2 > \varepsilon$)

$$\alpha_k = \frac{(A^T r_k, A^T r_k)}{(A^T p_k, A^T p_k)};$$

$$x_{k+1} = x_k + \alpha_k p_k; \quad r_{k+1} = r_k - \alpha_k A p_k;$$

$$\beta_k = \frac{(A^T r_{k+1}, A^T r_{k+1})}{(A^T r_k, A^T r_k)};$$

$$p_{k+1} = A^T r_{k+1} + \beta_k p_k;$$

$$k = k + 1;$$

end

算法 4.4 每次迭代需要计算三次矩阵与向量的乘法: Ap_k , $A^T p_k$ 和 $A^T r_{k+1}$. 注意到, CGNR 方法是极小化误差

$$\|e_k\|_{A^T A} = \|x_k - x^*\|_{A^T A},$$

式中: x^* 为精确解, 满足 $x^* = A^{-1}b$; $\|e_k\|_{A^T A}$ 为仿射空间

$$x_0 + \text{span}\{A^T r_0, (A^T A)A^T r_0, \dots, (A^T A)^{k-1}A^T r_0\} := x_0 + \mathcal{K}_k(A^T A, A^T r_0)$$

上关于残量 $b - Ax_k$ 的 2-范数. 因此, 可以比照 CG 方法的收敛性分析技巧对 CGNR 方法进行相应的收敛性分析.

2. CGNE 方法

对称不定或非对称非奇异的方程组 $Ax = b$ 还可转化为如下形式的对称正定的方程组, 即

$$AA^T y = b, \quad x = A^T y. \quad (4.50)$$

在“ y ”空间中, 对方程组 (4.50) 应用 CG 法, 得到如下形式:

选取初值 y_0 , 计算 $r_0 = b - AA^T y_0$; $p_0 = r_0$; $k = 0$;

while ($\|r_k\|_2 / \|r_0\| > \varepsilon$)

$$\alpha_k = \frac{(r_k, r_k)}{(p_k, AA^T p_k)} = \frac{(r_k, r_k)}{(A^T p_k, A^T p_k)};$$

$$y_{k+1} = y_k + \alpha_k p_k; \quad r_{k+1} = r_k - \alpha_k AA^T p_k;$$

$$\beta_k = \frac{(r_{k+1}, r_{k+1})}{(r_k, r_k)}; \quad p_{k+1} = r_{k+1} + \beta_k p_k;$$

$$k = k + 1;$$

end

作变量替换 $A^T y_k \rightarrow x_k$, $A^T p_k \rightarrow p_k$, 再化简上述有关式子则得到下面的关于法方程误差的共轭梯度法 (CGNE).

算法 4.4 (CGNE 方法) 给定初始向量 x_0 和容许误差限 $\varepsilon > 0$. 本算法计算 x_k , 使得 $\|r_k\|_2 / \|r_0\|_2 \leq \varepsilon$, 其中 $r_k = b - Ax_k$.

选取 x_0 ; 计算 $r_0 = b - Ax_0$; $p_0 = A^T r_0$; $k = 0$;

while ($\|r_k\|_2 / \|r_0\| > \varepsilon$)

$$\begin{aligned}\alpha_k &= \frac{(\mathbf{r}_k, \mathbf{r}_k)}{(\mathbf{p}_k, \mathbf{p}_k)}; \quad \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k; \\ \mathbf{r}_{k+1} &= \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k; \quad \beta_k = \frac{(\mathbf{r}_{k+1}, \mathbf{r}_{k+1})}{(\mathbf{r}_k, \mathbf{r}_k)}; \\ \mathbf{p}_{k+1} &= \mathbf{A}^T \mathbf{r}_{k+1} + \beta_k \mathbf{p}_k; \\ k &= k + 1;\end{aligned}$$

end

算法 4.4 每次迭代需要计算两次矩阵与向量的乘法: $\mathbf{A} \mathbf{p}_k$ 和 $\mathbf{A}^T \mathbf{r}_{k+1}$. 此外, CGNE 方法是极小化误差

$$\|\tilde{\mathbf{e}}_k\|_{\mathbf{A}^T \mathbf{A}} = \|\mathbf{y}_k - \mathbf{y}^*\|_{\mathbf{A}^T \mathbf{A}},$$

式中: \mathbf{y}^* 满足 $\mathbf{y}^* = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{b}$; $\|\tilde{\mathbf{e}}_k\|_{\mathbf{A}^T \mathbf{A}}$ 为仿射空间

$$\mathbf{x}_k \in \mathbf{x}_0 + \text{span}\{\mathbf{A}^T \mathbf{r}_0, (\mathbf{A}^T \mathbf{A}) \mathbf{A}^T \mathbf{r}_0, \dots, (\mathbf{A}^T \mathbf{A})^{k-1} \mathbf{A}^T \mathbf{r}_0\} \equiv \mathbf{x}_0 + \mathcal{K}_k(\mathbf{A}^T \mathbf{A}, \mathbf{A}^T \mathbf{r}_0)$$

上的 2-范数 $\|\mathbf{x}_k - \mathbf{x}^*\|_2$.

一般来说, CGNR 和 CGNE 这两种方法的困难或者说缺点在于矩阵的条件数变成了原来的平方, 然而, 在有些情况下它们是很有效的. 可参看文献 [25] 中的论述.

例 4.3 考虑方程组 $\mathbf{A} \mathbf{x} = \mathbf{b}$, 其中

$$\mathbf{A} = \begin{bmatrix} 12 & 3 & 2 & & & \\ -3 & 12 & 3 & 2 & & \\ -2 & -3 & 12 & 3 & 2 & \\ & \ddots & \ddots & \ddots & \ddots & \ddots \\ & & -2 & -3 & 12 & 3 \\ & & & -2 & -3 & 12 \end{bmatrix} \in \mathbb{R}^{n \times n}, \quad \mathbf{b} = \mathbf{A} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ \vdots \\ 1 \\ 1 \end{bmatrix} \in \mathbb{R}^n.$$

取 $n = 1000$, 分别将 CGNR 和 CGNE 方法应用到该线性方程组上, 两个迭代法均在 10 步后收敛 ($\varepsilon = 10^{-10}$), 计算得到的近似解 $\hat{\mathbf{x}}$ 和 $\tilde{\mathbf{x}}$ 与真解 \mathbf{x}^* 之间的绝对值误差分别为

$$\|\hat{\mathbf{x}} - \mathbf{x}^*\|_2 = 3.4705 \times 10^{-10}, \quad \|\tilde{\mathbf{x}} - \mathbf{x}^*\|_2 = 3.4515 \times 10^{-10}.$$

残量分别为

$$\|\mathbf{b} - \mathbf{A} \hat{\mathbf{x}}\|_2 = 4.5764 \times 10^{-9}, \quad \|\mathbf{b} - \mathbf{A} \tilde{\mathbf{x}}\|_2 = 4.6018 \times 10^{-9}.$$

迭代过程的收敛轨迹如图 4.3 所示, 其中横坐标为迭代步数 k , 纵坐标为相对残差 $\|\mathbf{r}_k\|_2 / \|\mathbf{r}_0\|_2$, 这里 \mathbf{r}_k 是第 k 步得到的残差向量.

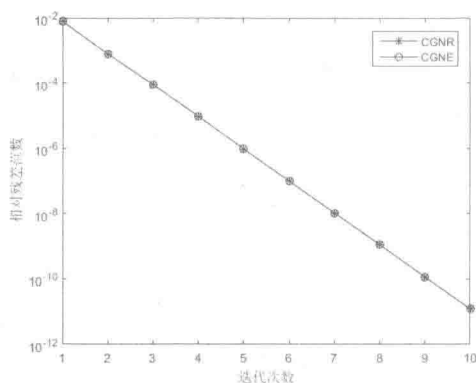


图 4.3 CGNR 和 CGNE 方法的收敛特性

4.2 广义极小残量法

广义极小残量法是目前求解大型稀疏非对称线性方程组的最常用方法之一, 在文献中常称为 GMRES 方法 (Generalized Minimal RESidual Method). 本节主要介绍这一算法的详细计算步骤及其收敛性理论.

4.2.1 GMRES 方法

考虑如下的线性方程组

$$Ax = b, \quad (4.51)$$

式中: 矩阵 $A \in \mathbb{R}^{n \times n}$ 和向量 $b \in \mathbb{R}^n$ 是给定的, 而 $x \in \mathbb{R}^n$ 是待求的未知向量. 这里假定系数矩阵 A 是非奇异的大型稀疏矩阵, 而且 $A^T \neq A$. 广义极小残量法是求 $x_k \in x_0 + \mathcal{K}_k(A, r_0)$, 使得

$$\|r_k\|_2 = \min\{\|b - Ax\|_2 : x \in x_0 + \mathcal{K}_k(A, r_0)\}, \quad (4.52)$$

式中: $r_k = b - Ax_k$, 即求 $x_k \in x_0 + \mathcal{K}_k(A, r_0)$, 使得残量 r_k 的 2-范数达到最小.

由于现在 A 是非对称的, 所以只能借助 A 的 Arnoldi 正交化过程极小化问题 (4.52). 回顾一下 Arnoldi 正交化过程 (见算法 2.11):

选取初始向量 v_1 , 使得 $\|v_1\|_2 = 1$;

for $j = 1 : k$

for $i = 1 : j$

$$h_{ij} = (Av_j, v_i);$$

$$\tilde{v}_{j+1} = Av_j - \sum_{i=1}^j h_{ij}v_i; \quad (4.53)$$

$$h_{j+1,j} = \|\tilde{v}_{j+1}\|_2; \quad v_{j+1} = \tilde{v}_{j+1}/h_{j+1,j};$$

end

end

不难发现, 式 (4.53) 可以写成

$$\mathbf{A}\mathbf{V}_k = \mathbf{V}_k\mathbf{H}_k + \beta_k\mathbf{v}_{k+1}\mathbf{e}_k^T = \mathbf{V}_{k+1}\widetilde{\mathbf{H}}_{k+1,k}, \quad (4.54)$$

式中: $\mathbf{V}_{k+1} = [\mathbf{V}_k, \mathbf{v}_{k+1}] \in \mathbb{R}^{n \times (k+1)}$ 满足 $\mathbf{V}_{k+1}^T \mathbf{V}_{k+1} = \mathbf{I}_{k+1}$; 矩阵

$$\widetilde{\mathbf{H}}_{k+1,k} = \begin{bmatrix} \mathbf{H}_k \\ \beta_k\mathbf{e}_k^T \end{bmatrix} \in \mathbb{R}^{(k+1) \times k}$$

为上 Hessenberg 矩阵, 其中 $\beta_k = h_{k+1,k}$,

$$\mathbf{H}_k = \begin{bmatrix} h_{11} & h_{12} & \cdots & h_{1k} \\ h_{21} & h_{22} & \cdots & h_{2k} \\ & h_{32} & \cdots & h_{3k} \\ & & \ddots & \vdots \\ & & & h_{k,k-1} & h_{k,k} \end{bmatrix}.$$

注意到 $\mathcal{K}_k(\mathbf{A}, \mathbf{r}_0) = \mathcal{R}(\mathbf{V}_k)$, 而且对任意的 $\mathbf{x} = \mathbf{x}_0 + \mathbf{V}_k\mathbf{z} \in \mathbf{x}_0 + \mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$, 有

$$\begin{aligned} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2 &= \|\mathbf{b} - \mathbf{A}\mathbf{x}_0 - \mathbf{A}\mathbf{V}_k\mathbf{z}\|_2 = \|\mathbf{r}_0 - \mathbf{A}\mathbf{V}_k\mathbf{z}\|_2 \\ &= \|\beta\mathbf{V}_{k+1}\mathbf{e}_1^{(k+1)} - \mathbf{V}_{k+1}\widetilde{\mathbf{H}}_{k+1,k}\mathbf{z}\|_2 \\ &= \|\beta\mathbf{e}_1^{(k+1)} - \widetilde{\mathbf{H}}_{k+1,k}\mathbf{z}\|_2, \end{aligned}$$

式中: $\beta = \|\mathbf{r}_0\|_2$; $\mathbf{e}_1^{(k+1)}$ 表示第 1 个分量为 1、其余分量均为 0 的 $k+1$ 维列向量. 由此可知, 极小化问题 (4.52) 等价于求 $\mathbf{z}_k \in \mathbb{R}^k$, 使得

$$\|\beta\mathbf{e}_1^{(k+1)} - \widetilde{\mathbf{H}}_{k+1,k}\mathbf{z}_k\|_2 = \min \{ \|\beta\mathbf{e}_1^{(k+1)} - \widetilde{\mathbf{H}}_{k+1,k}\mathbf{z}\|_2 : \mathbf{z} \in \mathbb{R}^k \}. \quad (4.55)$$

一旦这样的 \mathbf{z}_k 已经求得, 则所需的 \mathbf{x}_k 就为 $\mathbf{x}_k = \mathbf{x}_0 + \mathbf{V}_k\mathbf{z}_k$.

最小二乘问题 (4.55) 也可以用 $\widetilde{\mathbf{H}}_{k+1,k}$ 的 QR 分解来求解. 由于 $\widetilde{\mathbf{H}}_{k+1,k}$ 是上 Hessenberg 矩阵, 可以计算 k 个 Givens 变换

$$\mathbf{G}_i = \text{diag} \left(\mathbf{I}_{i-1}, \begin{bmatrix} c_i & s_i \\ -s_i & c_i \end{bmatrix}, \mathbf{I}_{k-i} \right), \quad c_i^2 + s_i^2 = 1,$$

使得

$$(\mathbf{G}_k\mathbf{G}_{k-1}\cdots\mathbf{G}_2\mathbf{G}_1)\widetilde{\mathbf{H}}_{k+1,k} = \begin{bmatrix} \mathbf{R}_k \\ \mathbf{0} \end{bmatrix}, \quad (4.56)$$

式中: \mathbf{R}_k 为非奇异的上三角矩阵 (因为 $\widetilde{\mathbf{H}}_{k+1,k}$ 的次对角元均不为零). 令

$$\mathbf{G} = \mathbf{G}_k\mathbf{G}_{k-1}\cdots\mathbf{G}_2\mathbf{G}_1, \quad \begin{bmatrix} \mathbf{t}_k \\ \rho_k \end{bmatrix} = \mathbf{G}(\beta\mathbf{e}_1), \quad \mathbf{t}_k = (\tau_1, \tau_2, \cdots, \tau_k)^T, \quad (4.57)$$

则 G 是 $k+1$ 阶正交矩阵, 而且直接计算, 有

$$\begin{cases} \tau_1 = \beta c_1, \\ \tau_i = (-1)^{i-1} \beta s_1 s_2 \cdots s_{i-1} c_i, \quad i = 2, 3, \cdots, k, \\ \rho_k = (-1)^k \beta s_1 s_2 \cdots s_k. \end{cases} \quad (4.58)$$

由此立即可得最小二乘问题 (4.55) 的解为

$$z_k = R_k^{-1} t_k. \quad (4.59)$$

此外, 此时的残量范数为

$$\|r_k\|_2 = \|\beta e_1^{(k+1)} - \widetilde{H}_{k+1,k} z_k\|_2 = |\rho_k|. \quad (4.60)$$

综合上面的讨论, 可以将广义极小残量法 (GMRES) 的步骤总结如下:

- (1) 令 $v_1 = r_0 / \|r_0\|_2$, 产生一个长度为 k 的 Arnoldi 分解 (4.54).
- (2) 利用 Givens 变换求 $\widetilde{H}_{k+1,k}$ 的 QR 分解 (4.56), 并按式 (4.58) 求得向量 t_k 和数 ρ_k .
- (3) 用回代法求解上三角方程组 $R_k z_k = t_k$, 得 z_k .
- (4) 计算 $x_k = x_0 + V_k z_k$.
- (5) 若 $|\rho_k|/\beta < \varepsilon$ (事先给定的误差界), 则终止; 否则增加 k 的值, 重复上面的过程.

上述 GMRES 方法可写成如下便于编程的算法形式.

算法 4.5 (GMRES 方法) 给定矩阵 $A \in \mathbb{R}^{n \times n}$, 向量 $b \in \mathbb{R}^n$ 和允许误差 $\varepsilon > 0$, 取初始向量 x_0 . 本算法计算 $x_k \in \mathbb{R}^n$, 使得 $\|r_k\|_2 / \|r_0\|_2 \leq \varepsilon$, 其中 $r_k = b - Ax_k$.

步 1, 计算初始残量 $r_0 = b - Ax_0$, $\beta = \|r_0\|_2$ 和初始正交化向量 $v_1 = r_0/\beta$, $\xi = e_1 = (1, 0, \cdots, 0)^T$. 置 $k := 1$.

步 2, 用 Arnoldi 过程计算 v_{k+1} 和 $h_{i,k}$ ($i = 1, 2, \cdots, k+1$). 将 Givens 变换 G_i 作用于矩阵 $\widetilde{H}_{i+1,i}$ 的最后一列:

$$\begin{bmatrix} h_{i,k} \\ h_{i+1,k} \end{bmatrix} := \begin{bmatrix} c_i & s_i \\ -s_i & c_i \end{bmatrix} \begin{bmatrix} h_{i,k} \\ h_{i+1,k} \end{bmatrix}, \quad i = 1, \cdots, k-1.$$

步 3, 计算第 k 次 Givens 变换 G_k 中的 c_k 和 s_k :

$$\begin{aligned} c_k &= \frac{h_{k,k}}{\sqrt{h_{k,k}^2 + h_{k+1,k}^2}}, \quad s_k = \frac{h_{k+1,k}}{\sqrt{h_{k,k}^2 + h_{k+1,k}^2}}, \\ \left(\tau_k &= \frac{h_{k+1,k}}{h_{k,k}}, \quad c_k = \frac{1}{\sqrt{1 + \tau_k^2}}, \quad s_k = c_k \tau_k \right). \end{aligned} \quad (4.61)$$

步 4, 对 $k+1$ 维向量 ξ 的最后两个元素和矩阵 $\widetilde{H}_{k+1,k}$ 分别作用第 k 次 Givens 变换 G_k , 有

$$\begin{bmatrix} \xi_k \\ \xi_{k+1} \end{bmatrix} := \begin{bmatrix} c_k & s_k \\ -s_k & c_k \end{bmatrix} \begin{bmatrix} \xi_k \\ 0 \end{bmatrix} = \begin{bmatrix} c_k \xi_k \\ -s_k \xi_k \end{bmatrix},$$

$$\begin{bmatrix} h_{k,k} \\ h_{k+1,k} \end{bmatrix} := \begin{bmatrix} c_k & s_k \\ -s_k & c_k \end{bmatrix} \begin{bmatrix} h_{k,k} \\ h_{k+1,k} \end{bmatrix} = \begin{bmatrix} c_k h_{k,k} + s_k h_{k+1,k} \\ 0 \end{bmatrix}. \quad (4.62)$$

步 5, 若 $|\rho_k|/\beta = |\xi_{k+1}|/\beta \leq \varepsilon$, 求解关于 z_k 的上三角方程组 $H_{k,k} z_k = [\beta \xi]_{k \times 1}$, 计算近似解向量 $x_k = x_0 + V_k z_k$, 停算; 否则, 置 $k := k+1$, 转步 2.

算法 4.5 的 MATLAB 程序如下:

```
%GMRES方法程序-mgmres.m
function [x,k,time,res,resvec,flag]=mgmres(A,b,x,max_it,tol)
tic; flag=0; r=b-A*x; beta=norm(r); %计算残差
n=length(b); e1=zeros(n,1); e1(1)=1.0;
res=norm(r)/beta; resvec(1)=res;
V(:,1)=r/beta; xi=beta*e1; k=0;
while (k<=max_it)
    k=k+1; w=A*V(:,k);
    for i=1:k %修正Arnoldi过程
        H(i,k)=w'*V(:,i);
        w=w-H(i,k)*V(:,i);
    end
    H(k+1,k)=norm(w);
    if abs(H(k+1,k))/beta<tol,
        return;
    else
        V(:,k+1)=w/H(k+1,k);
    end
    for i=1:k-1,
        temp=c(i)*H(i,k)+s(i)*H(i+1,k);
        H(i+1,k)=-s(i)*H(i,k)+c(i)*H(i+1,k);
        H(i,k)=temp;
    end
    [c(k),s(k),H(k,k)]=givens(H(k,k), H(k+1,k)); %第k次Givens变换
    xi(k+1)=-s(k)*xi(k);
    xi(k)=c(k)*xi(k); H(k+1,k)=0.0;
    res=abs(xi(k+1))/beta;
    resvec(k+1)=res;
```

```

if (res<=tol ),
    y=H(1:k,1:k)\xi(1:k); x=x+V(:,1:k)*y;
    break; %跳出循环
end
end
if (res>tol), flag=1; end; %不收敛
time=toc;

```

例 4.4 考虑系数矩阵 A 和右端项 b 分别为

$$A = \begin{bmatrix} 1 & -1 & 0 & \cdots & 0 & n \\ 1 & 2 & -1 & \ddots & & 0 \\ 0 & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & 0 \\ 0 & & \ddots & 1 & n-1 & -1 \\ -n & 0 & \cdots & 0 & 1 & n \end{bmatrix}, \quad b = A \begin{bmatrix} 1 \\ 1 \\ \vdots \\ \vdots \\ 1 \\ 1 \end{bmatrix}$$

的线性方程组 (4.51). 显然, 该方程组的真解为 $x^* = (1, 1, \dots, 1)^T$. 取 $n = 10^3$, 应用 GMRES 算法到该方程组上, 迭代 172 步后得到的近似解 $\tilde{x} = x_{172}$ 满足

$$\|\tilde{x} - x^*\|_2 = 1.1427 \times 10^{-7}.$$

迭代过程的收敛轨迹如图 4.4 所示, 其中横坐标为迭代步数 k , 纵坐标为相对残差 $\|r_k\|_2 / \|r_0\|_2$, 这里 r_k 是第 k 步得到的残差向量.

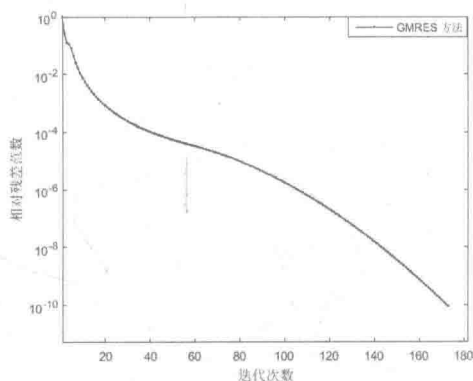


图 4.4 GMRES 方法的收敛特性

实际使用 GMRES 方法时的主要问题是 k 不能太大, 这是因为这一方法的内存需求量为 $O(kn)$, 而计算机的内存又是有限的, 这样就会出现近似解 x_k 还不满足精

度要求, 而 k 已经不能再增加的情形. 解决这一问题的一个简单而行之有效的办法就是重新开始技术. 它的基本思想是, 先选定一个不太大的正整数 m , 用 GMRES 方法产生 x_m , 然后再以 x_m 作为初始向量重新开始. 这就是重新开始 GMRES 方法, 通常记为 GMRES(m) 方法, 具体计算过程可简述如下.

算法 4.6 (GMRES(m) 方法) 给定矩阵 $A \in \mathbb{R}^{n \times n}$, 向量 $b \in \mathbb{R}^n$, 正整数 m , 初始向量 x_0 和允许误差 $\varepsilon > 0$. 本算法计算 $x_m \in \mathbb{R}^n$, 使得 $\|b - Ax_m\|_2 / \|b - Ax_0\|_2 \leq \varepsilon$.

选取 x_0 ; 计算 $r_0 = b - Ax_0$; $\beta = \|r_0\|_2$; $\rho_m := \beta$; $x_m := x_0$;

while ($\rho_m / \beta > \varepsilon$)

$r_0 = b - Ax_m$; $\beta_0 = \|r_0\|_2$; $v_1 = r_0 / \beta_0$;

以 v_1 为初始向量产生一个长度为 m 的 Arnoldi 分解

$$AV_m = V_{m+1} \widetilde{H}_{m+1,m};$$

计算 $\widetilde{H}_{m+1,m}$ 的 QR 分解: $\widetilde{H}_{m+1,m} = G^T \begin{bmatrix} R_m \\ 0 \end{bmatrix}$;

按式 (4.58) 计算 t_m 和 ρ_m ;

求解 $R_m z_m = t_m$ 得到 z_m ;

计算 $x_m = x_0 + V_m z_m$;

end

这一算法是目前求解大型稀疏非对称线性方程组的常用方法之一. 至于算法 4.6 中的 m 取多大为好, 现在还没有理论上的结果. 在理论上仅可以保证, 在系数矩阵 A 有正定的对称部分时, 对任意的 m , 算法 4.6 总是收敛的 (见后面的收敛性定理). 但对于一般情形, 并非对任意的 m 总能保证其收敛. 例如, 对线性方程组

$$Ax := \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} := b,$$

利用 GMRES(1) 求解, 不论循环多少次, 总有 $x_k = 0$, 而 $\|Ax_k - b\|_2 = \sqrt{2}$, 永远不会得到原方程组之真解满足要求的近似解. 但如果用 GMRES 求解, 则只需两步就可以得到该方程组的精确解.

GMRES(m) 的 MATLAB 程序如下:

```
%重新开始GMRES方法-GMRES(m)-gmres.m
function [x,out,int,time,res,resvec,flag]=gmres(A,b,x,m,max_it,tol)
tic; flag=0; int=0; r= b-A*x; %计算残差
beta=norm(r); res=norm(r)/beta; resvec(1)=res;
n=length(b); %m=resttrt;
V(1:n,1:m+1)=zeros(n,m+1);
H(1:m+1,1:m)=zeros(m+1,m);
e1=zeros(n,1); e1(1)=1.0;
```

```

c(1:m)=zeros(m,1); s(1:m)=zeros(m,1);
for k=1:max_it
    V(:,1)=r/norm(r);
    xi=norm(r)*e1;
    for j=1:m %用Arnoldi方法构造正交基
        w=A*V(:,j);
        for i=1:j
            H(i,j)=w'*V(:,i);
            w=w-H(i,j)*V(:,i);
        end
        H(j+1,j)=norm(w);
        if abs(H(j+1,j))/beta<tol,
            return;
        else
            V(:,j+1)=w/H(j+1,j);
        end
        for i=1:j-1 %第i次Givens变换
            temp=c(i)*H(i,j)+s(i)*H(i+1,j);
            H(i+1,j)=-s(i)*H(i,j)+c(i)*H(i+1,j);
            H(i,j)=temp;
        end
        [c(j),s(j),H(j,j)]=givens(H(j,j),H(j+1,j));%第j次Givens变换
        xi(j+1)=-s(j)*xi(j);xi(j)=c(j)*xi(j);
        H(j+1,j)=0.0;res=abs(xi(j+1))/beta;
        resvec((k-1)*m+j+1)=res;
        if(res<=tol)
            y=H(1:j,1:j)\xi(1:j);
            x=x+V(:,1:j)*y;
            break; %跳出内循环
        end
    end
    if (res<tol )
        out=k; int=j; break; %跳出外循环
    end
    y=H(1:m,1:m)\xi(1:m);
    x=x+V(:,1:m)*y;
    r=b-A*x ;
end
if (res>tol), flag=1; end; %不收敛
time=toc;

```

例 4.5 用 GMRES(m) 方法求解例 4.4 的线性方程组. 表 4.1 列出了对于不同的 m , 达到收敛 (容许误差为 $\varepsilon = 10^{-10}$) 所需要的外迭代次数、内迭代次数、总迭代次数和 CPU 时间 (s). 从表 4.1 中的数据可看出, 对于本例而言, 与 GMRES 方法相比, 对

表 4.1 GMRES(m) 方法对 m 的依赖性

m	外迭代次数	内迭代次数	总迭代次数	CPU 时间	相对残差
10	47	3	463	0.2621	9.8273e-11
20	14	12	272	0.1488	9.1166e-11
30	9	8	248	0.1327	9.3534e-11
40	6	27	227	0.1267	9.4923e-11
50	5	19	219	0.1244	9.9472e-11
60	4	26	206	0.1203	9.9062e-11
GMRES	1	172	172	0.1651	8.8473e-11

于比较小的 m 值, GMRES(m) 方法并不占优势. 事实上, 前面已经指出, 只有在无法使用 GMRES 方法时 (如超出计算机的内存), 才考虑使用 GMRES(m) 方法, 并且 m 不宜取得太小.

此外, 对于不同的 m 值, 迭代过程的收敛轨迹如图 4.5 所示, 其中横坐标为迭代步数 k , 纵坐标为相对残差 $\|r_k\|_2/\|r_0\|_2$, 这里 r_k 是第 k 步得到的残差向量.

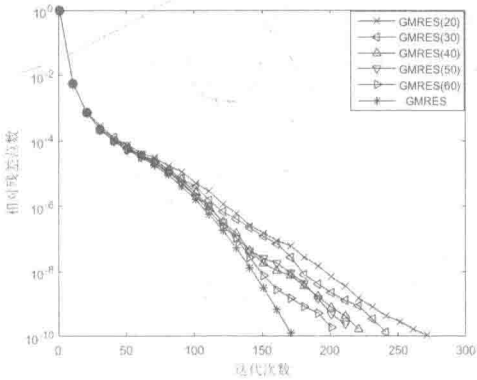


图 4.5 GMRES 方法的收敛特性

4.2.2 预处理 GMRES 方法

回顾共轭梯度法求解方程组 $Ax = b$ 时预处理技术的重要性, 考虑预处理广义极小残量法 (简记为 PGMRES 方法). 因为共轭梯度法适用于对称正定矩阵 A , 所以预处理矩阵 M 需要保持对称. 对于非对称矩阵 A , 预处理矩阵 M 就没有必要是对称的. 此时方程组 $Ax = b$ 被预处理为

$$M^{-1}Ax = M^{-1}b. \tag{4.63}$$

这里 M 应该取为 A 的一个近似矩阵, 以保证用 GMRES 方法求解预处理方程组 (4.63) 比求解原方程组 $Ax = b$ 有更快的收敛速度.

由于 GMRES 方法只涉及到矩阵与向量的乘法, 所以为了保证 M^{-1} 与某一个向量 r 的乘法 $z = M^{-1}r$ 容易计算, 即 $Mz = r$ 容易求解. 这启发选取 M 为 (块) 对角矩阵或 (块) 三角形矩阵, 如 Jacobi 预处理矩阵

$$M = D = \text{diag}\{a_{11}, \dots, a_{nn}\},$$

Gauss-Seidel 预处理矩阵

$$M = D - L$$

以及 SOR 预处理矩阵

$$M = \frac{1}{\omega}(D - \omega L),$$

式中:

$$L = - \begin{bmatrix} 0 & & & \\ a_{21} & 0 & & \\ \vdots & \ddots & \ddots & \\ a_{n1} & \cdots & a_{n,n-1} & 0 \end{bmatrix}.$$

对方程组 (4.63) 用 GMRES 方法, 得到如下 PGMRES 方法.

算法 4.7 (PGMRES 方法) 给定矩阵 $A \in \mathbb{R}^{n \times n}$, 向量 $b \in \mathbb{R}^n$, 预处理子 M , 初始向量 x_0 和允许误差 $\varepsilon > 0$. 本算法计算 $x_k \in \mathbb{R}^n$, 使得 $\|r_k\|_2 / \|r_0\|_2 \leq \varepsilon$, 其中 $r_k = b - Ax_k$.

$$r_0 = M^{-1}(b - Ax_0); \beta = \|r_0\|_2;$$

for $k = 1, 2, \dots$

对 $K_k(M^{-1}A, v_1)$ 用 Arnoldi 分解计算 $V_k = [v_1, v_2, \dots, v_k]$

以及 $(k+1) \times k$ 阶矩阵 $\widetilde{H}_{k+1,k}$.

计算 $\widetilde{H}_{k+1,k}$ 的 QR 分解: $\widetilde{H}_{k+1,k} = G^T \begin{bmatrix} R_k \\ 0 \end{bmatrix};$

按式 (4.58) 计算 t_k 和 ρ_k ;

if $|\rho_k|/\beta \leq \varepsilon$

求解 $R_k z_k = t_k$ 得到 z_k ;

计算 $x_k = x_0 + V_k z_k$; 停算.

end

end

对方程组 (4.63) 用 GMRES(m) 得到如下算法:

算法 4.8 (PGMRES(m) 方法) 给定矩阵 $A \in \mathbb{R}^{n \times n}$, 向量 $b \in \mathbb{R}^n$, 正整数 m , 预处理子 M 和允许误差 $\varepsilon > 0$, 选取初始向量 $x_0 \in \mathbb{R}^n$. 本算法计算 $x_k \in \mathbb{R}^n$, 使得 $\|r_k\|_2 / \|r_0\|_2 \leq \varepsilon$, 其中 $r_k = b - Ax_k$.

步 1, 计算残量 $r_0 = M^{-1}(b - Ax_0)$, $\beta = \|r_0\|_2$ 和初始正交化向量 $v_1 = r_0/\beta$.

步 2, 对 $K_m(M^{-1}A, v_1)$ 用 Arnoldi 分解计算 $V_m = [v_1, v_2, \dots, v_m]$ 以及 $(m+1) \times m$ 阶矩阵 $\widetilde{H}_{m+1,m}$.

步 3, 计算满足极小化问题

$$\min \{ \|\beta e_1^{(m+1)} - \widetilde{H}_{m+1,m} z\|_2 : z \in \mathbb{R}^m \}$$

的最小二乘解 z_m , 令 $x_m = x_0 + V_m z_m$. 若 x_m 达到精度要求, 停算.

步 4, 置 $x_0 := x_m$, 转步 1.

下面通过数值例子, 观察 PGMRES 方法的效果.

例 4.6 仍考虑例 4.4 的线性方程组. 取预处理矩阵 $M = \text{diag}(A)$, 应用 PGMRES 方法到该方程组上, 迭代 12 步后得到的近似解 $\tilde{x} = x_{12}$ 满足

$$\|\tilde{x} - x_*\|_2 = 1.7407 \times 10^{-8}.$$

迭代过程的收敛轨迹如图 4.6 所示, 其中横坐标为迭代步数 k , 纵坐标为相对残差 $\|r_k\|_2 / \|r_0\|_2$, 这里 r_k 是第 k 步得到的残差向量.

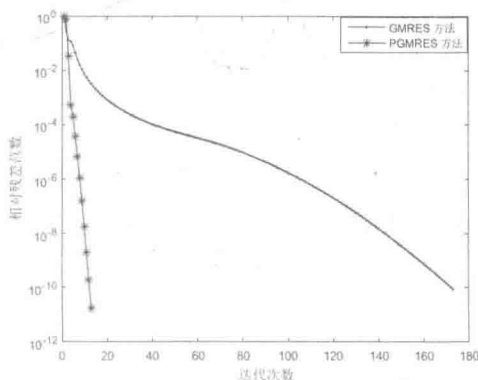


图 4.6 GMRES 和 PGMRES 的收敛特性

例 4.7 用 PGMRES(m) 方法求解例 4.4 的线性方程组, 取预处理矩阵 $M = \text{diag}(A)$. 表 4.2 列出了对于不同的 m , 达到收敛 (容许误差为 $\varepsilon = 10^{-10}$) 所需要的外迭代次数、内迭代次数、总迭代次数和 CPU 时间 (s).

此外, 对于不同的 m 值, 迭代过程的收敛轨迹如图 4.7 所示, 其中横坐标为迭代步数 k , 纵坐标为相对残差 $\|r_k\|_2 / \|r_0\|_2$, 这里 r_k 是第 k 步得到的残差向量.

表 4.2 PGMRES(*m*) 方法对 *m* 的依赖性

<i>m</i>	外迭代次数	内迭代次数	总迭代次数	CPU 时间	相对残差
3	21	3	63	0.1692	7.3458e-11
6	4	3	21	0.0449	3.8001e-11
9	2	5	14	0.0290	8.0750e-11
12	1	12	12	0.0238	1.7551e-11
PGMRES	1	12	12	0.0394	1.7551e-11

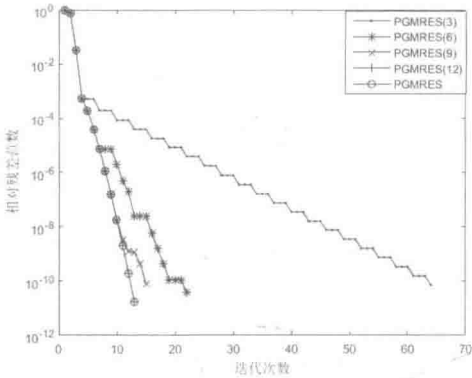


图 4.7 PGMRES(*m*) 方法对 *m* 的依赖性

4.2.3 收敛性分析

前一节 CG 方法的收敛性都是关于对称矩阵的. 当系数矩阵非对称时, 对 Krylov 子空间方法 (比如 GMRES 方法) 的收敛性进行分析就会遇到很大的困难. 首先来看 GMRES 方法的一个重要性质.

定理 4.4 GMRES 方法不会发生中断.

证明 若 $h_{k+1,k} \neq 0$, 则计算过程直至第 k 步都不会中断. 事实上, 当 $h_{k+1,k} \neq 0$ 时, 由式 (4.61) 和式 (4.62), \mathbf{R}_k 的对角元满足

$$r_{k,k} = h_{k,k} := c_k h_{k,k} + s_k h_{k+1,k} = \sqrt{h_{k,k}^2 + h_{k+1,k}^2} > 0.$$

故正交化可进行, 极小化问题可解.

由此可知, 只有当 $h_{k+1,k} = 0$ 时, 计算过程才在第 k 步中断, 此时向量 \mathbf{v}_{k+1} 不能构造. 但此时成立

$$\mathbf{A}\mathbf{V}_k = \mathbf{V}_k\mathbf{H}_k,$$

故 $\sigma(\mathbf{H}_k) \subset \sigma(\mathbf{A})$. 由于 \mathbf{A} 非奇异, 故 \mathbf{H}_k 也非奇异. 在第 k 步, 极小化问题 (4.55) 变为

$$\|\beta \mathbf{v}_1 - \mathbf{A}\mathbf{V}_k \mathbf{z}\|_2 = \|\beta \mathbf{v}_1 - \mathbf{V}_k \mathbf{H}_k \mathbf{z}\|_2 = \|\mathbf{V}_k (\beta \mathbf{e}_1^{(k)} - \mathbf{H}_k \mathbf{z})\|_2$$

$$= \|\beta e_1^{(k)} - H_k z\|_2.$$

而 H_k 非奇异, 当 $z_k = H_k^{-1}(\beta e_1^{(k)})$ 时, $\|\beta v_1 - A V_k z\|_2$ 达到极小值 0, 即 $\|r_k\|_2 = 0$, 故

$$x_k = x_0 + V_k z_k$$

是精确解, 换言之, GMRES 方法若在第 k 步中断, 则在这一步已得到精确解. 证毕. \square

下面的定理是 GMRES 方法的另一个性质.

定理 4.5 设 $\{x_i\}$ 是 GMRES 方法产生的迭代序列. 若 x_k 是精确解, 而 $x_i (i < k)$ 不是精确解, 则算法在第 k 步中断.

证明 由假设, 有

$$r_i \neq 0, \quad i = 1, 2, \dots, k-1, \quad \text{而} \quad r_k = 0.$$

但 $\|r_k\|_2$ 是 $\tilde{\xi}_k = \beta V_k e_1^{(k+1)}$ 最后一个分量的绝对值, 即它是

$$\begin{bmatrix} 1 & & & & \\ & \ddots & & & \\ & & 1 & & \\ & & & c_k & s_k \\ & & & -s_k & c_k \end{bmatrix} \begin{bmatrix} \tilde{\xi}_{k-1} \\ \\ \\ 0 \end{bmatrix}$$

的第 $k+1$ 个分量的绝对值, 即

$$\|r_k\|_2 = |s_k e_k^T \tilde{\xi}_{k-1}|. \quad (4.64)$$

但

$$|e_k^T \tilde{\xi}_{k-1}| = \|r_{k-1}\|_2 \neq 0, \quad (4.65)$$

故 $s_k = 0$. 由于

$$s_k = \frac{h_{k+1,k}}{\sqrt{h_{k,k}^2 + h_{k+1,k}^2}},$$

故有 $h_{k+1,k} = 0$, 这样 $\tilde{v}_{k+1} = 0$ (由式 (4.53) 所定义), 算法中断. 证毕. \square

由式 (4.64) 和式 (4.65) 容易得到下面的推论.

推论 4.2 GMRES 方法的残量范数 $\|r_k\|_2$ 有表达式

$$\|r_k\|_2 = \left(\prod_{i=1}^k |s_i| \right) \|r_0\|_2. \quad (4.66)$$

定理 4.6 初始残量 r_0 的最小多项式次数是 k 的充分必要条件是 $\tilde{v}_{k+1} = 0$ 且 $\tilde{v}_i \neq 0 (i = 1, 2, \dots, k)$, 其中 $\tilde{v}_{i+1} (i = 1, 2, \dots, k)$ 由式 (4.53) 所定义.

证明 若 $\mathbf{r}_0 = \|\mathbf{r}_0\|_2 \mathbf{v}_1$ 的最小多项式次数为 k , 则存在一个 k 次多项式 ϕ_k , 使得 $\phi_k(\mathbf{A})\mathbf{v}_1 = \mathbf{0}$, 且 ϕ_k 是次数最低者. 因此

$$\mathcal{K}_{k+1} = \text{span}\{\mathbf{v}_1, \mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}^k \mathbf{v}_1\} = \mathcal{K}_k.$$

由于 $\tilde{\mathbf{v}}_{k+1} \in \mathcal{K}_{k+1} = \mathcal{K}_k$, 且 $\tilde{\mathbf{v}}_{k+1} \perp \mathcal{K}_k$, 故 $\tilde{\mathbf{v}}_{k+1} = \mathbf{0}$. 另外, 若存在某个 $i (1 \leq i \leq k)$, 使 $\tilde{\mathbf{v}}_i = \mathbf{0}$, 则存在一个 $i-1$ 次多项式 ϕ_{i-1} 使 $\phi_{i-1}(\mathbf{A})\mathbf{v}_1 = \mathbf{0}$, 这与 \mathbf{v}_1 的最小多项式次数为 k 矛盾, 故 $\tilde{\mathbf{v}}_i \neq \mathbf{0}, 1 \leq i \leq k$.

反之, 设 $\tilde{\mathbf{v}}_{k+1} = \mathbf{0}$ 且 $\tilde{\mathbf{v}}_i \neq \mathbf{0}, i = 1, 2, \dots, k$. 由 $\tilde{\mathbf{v}}_{k+1} = \mathbf{0}$ 可知, 存在一个 k 次多项式 ϕ_k , 使 $\phi_k(\mathbf{A})\mathbf{v}_1 = \mathbf{0}$, 而且 k 是次数最小者. 否则, 当存在 $\phi_i (i < k)$ 使 $\phi_i(\mathbf{A})\mathbf{v}_1 = \mathbf{0}$ 时, 就有 $\mathcal{K}_{i+1} = \mathcal{K}_i$, 因此 $\tilde{\mathbf{v}}_{i+1} = \mathbf{0}$, 这与 $\tilde{\mathbf{v}}_i \neq \mathbf{0} (1 \leq i \leq k)$ 矛盾. 证毕. \square

进一步, 由上述三个定理可以推出:

推论 4.3 GMRES 方法在第 k 步产生的解 \mathbf{x}_k 是精确解的充分必要条件为下列诸等价条件:

- (1) 算法在第 k 步中断.
- (2) $h_{k+1,k} = 0$.
- (3) $\tilde{\mathbf{v}}_{k+1} = \mathbf{0}$.
- (4) \mathbf{r}_0 的最小多项式次数为 k .

推论 4.4 GMRES 方法得到的残量模序列 $\{\|\mathbf{r}_k\|_2\}$ 是单调下降的, 对于 n 阶方程组至多迭代 n 步即可得到精确解.

现在考虑 GMRES(m) 方法, 虽然它总能进行下去, 但可能不收敛. 当然, 当 m 充分大时是收敛的. 特别地, 当 $m = n$ 时, 一个重新开始迭代即可得到精确解.

1. GMRES(m) 的收敛性定理

现在考虑 GMRES(m) 方法的误差估计. 利用 Krylov 子空间的性质, 注意到任意的 $\mathbf{x} = \mathbf{x}_0 + \varphi_{k-1}(\mathbf{A})\mathbf{r}_0$ (即 $\mathbf{x} \in \mathbf{x}_0 + \mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$), 可导出

$$\begin{aligned} \|\mathbf{r}_k\|_2 &= \min \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2 = \min_{\varphi_{k-1} \in \mathcal{P}_{k-1}} \|\mathbf{b} - \mathbf{A}(\mathbf{x}_0 + \varphi_{k-1}(\mathbf{A})\mathbf{r}_0)\|_2 \\ &= \min_{\varphi_{k-1} \in \mathcal{P}_{k-1}} \|(I - \mathbf{A}\varphi_{k-1}(\mathbf{A}))\mathbf{r}_0\|_2 = \min_{\psi \in \mathcal{P}_k^0} \|\psi(\mathbf{A})\mathbf{r}_0\|_2, \end{aligned} \quad (4.67)$$

式中: \mathcal{P}_k^0 如式 (4.41) 中所定义.

定理 4.7 设 \mathbf{A} 是正定的 (即 \mathbf{A} 的对称部分是对称正定的), $m > 0$ 是任意给定的整数, 并假定在 GMRES(m) 中重新开始了 ℓ 次产生了近似解 $\mathbf{x}_m^{(\ell)}$, 则有

$$\frac{\|\mathbf{r}_m^{(\ell)}\|_2}{\|\mathbf{r}_0\|_2} \leq \left[\sqrt{1 - \frac{\lambda_{\min}^2(M)}{\lambda_{\max}(\mathbf{A}^T \mathbf{A})}} \right]^{\ell m}, \quad (4.68)$$

式中: $M = \frac{1}{2}(\mathbf{A} + \mathbf{A}^T)$, 其他的记号如前所述.

证明 因为 \mathbf{A} 是正定的, 即其对称部分

$$\mathbf{M} = \frac{1}{2}(\mathbf{A} + \mathbf{A}^T) \quad (4.69)$$

是对称正定的, 定义

$$\psi_\alpha(t) = (1 + \alpha t)^m, \quad \alpha \in \mathbb{R}. \quad (4.70)$$

显然有 $\psi_\alpha(t) \in \mathcal{P}_m^0$. 下面先给出 $\min_\alpha \|\psi_\alpha(\mathbf{A})\|_2$ 的上界估计.

现在任取一个 $\mathbf{u} \in \mathbb{R}^n$, $\|\mathbf{u}\|_2 = 1$ (即 \mathbf{u} 是单位向量), 并记 $\hat{\psi}_\alpha(t) = 1 + \alpha t$, 则有

$$\begin{aligned} \|\hat{\psi}_\alpha(\mathbf{A})\mathbf{u}\|_2^2 &= \mathbf{u}^T (\mathbf{I} + \alpha \mathbf{A})^T (\mathbf{I} + \alpha \mathbf{A}) \mathbf{u} \\ &= 1 + 2\alpha \mathbf{u}^T \mathbf{M} \mathbf{u} + \alpha^2 \mathbf{u}^T \mathbf{A}^T \mathbf{A} \mathbf{u}, \end{aligned}$$

从而, 当 $\alpha \geq 0$ 时, 有

$$\|\hat{\psi}_\alpha(\mathbf{A})\|_2 \geq \|\hat{\psi}_\alpha(\mathbf{A})\mathbf{u}\|_2 \geq 1, \quad (4.71)$$

而当 $\alpha \leq 0$ 时, 有

$$\|\hat{\psi}_\alpha(\mathbf{A})\mathbf{u}\|_2^2 \leq 1 + 2\alpha \lambda_{\min}(\mathbf{M}) + \alpha^2 \lambda_{\max}(\mathbf{A}^T \mathbf{A}), \quad (4.72)$$

式中: $\lambda_{\min}(\mathbf{M})$ 和 $\lambda_{\max}(\mathbf{A}^T \mathbf{A})$ 分别为 \mathbf{M} 的最小特征值和 $\mathbf{A}^T \mathbf{A}$ 的最大特征值.

注意到单位向量 \mathbf{u} 的任意性, 式 (4.72) 蕴涵着当 $\alpha \leq 0$ 时, 有

$$\|\hat{\psi}_\alpha(\mathbf{A})\|_2^2 \leq 1 + 2\alpha \lambda_{\min}(\mathbf{M}) + \alpha^2 \lambda_{\max}(\mathbf{A}^T \mathbf{A}). \quad (4.73)$$

不等式 (4.73) 的右边是关于 α 的二次函数, 在

$$\alpha = -\frac{\lambda_{\min}(\mathbf{M})}{\lambda_{\max}(\mathbf{A}^T \mathbf{A})} \leq 0$$

处达到最小值 $1 - \lambda_{\min}^2(\mathbf{M})/\lambda_{\max}(\mathbf{A}^T \mathbf{A})$, 从而有

$$\min_{\alpha < 0} \|\hat{\psi}_\alpha(\mathbf{A})\|_2^2 \leq 1 - \frac{\lambda_{\min}^2(\mathbf{M})}{\lambda_{\max}(\mathbf{A}^T \mathbf{A})} < 1.$$

再注意到式 (4.71), 有

$$\min_{\alpha \in \mathbb{R}} \|\hat{\psi}_\alpha(\mathbf{A})\|_2 \leq \sqrt{1 - \frac{\lambda_{\min}^2(\mathbf{M})}{\lambda_{\max}(\mathbf{A}^T \mathbf{A})}} \equiv \kappa. \quad (4.74)$$

如果从 \mathbf{x}_0 出发, 应用 GMRES 迭代 m 步得到 \mathbf{x}_m , 则对 $k = m$ 应用式 (4.67), 有

$$\begin{aligned} \|\mathbf{r}_m\|_2 &\leq \min_\alpha \|\psi_\alpha(\mathbf{A})\mathbf{r}_0\|_2 \leq \min_\alpha \|\psi_\alpha(\mathbf{A})\|_2 \|\mathbf{r}_0\|_2 \\ &\leq \min_\alpha \|\hat{\psi}_\alpha(\mathbf{A})\|_2^m \|\mathbf{r}_0\|_2 \leq \kappa^m \|\mathbf{r}_0\|_2, \end{aligned} \quad (4.75)$$

其中最后一个不等式用到了式 (4.74).

假定以 $x_m^{(1)} = x_m$ 再重新开始, 用 GMRES 迭代产生 $x_m^{(2)}$, 则由式 (4.75) 又有

$$\|r_m^{(2)}\|_2 \leq \kappa^m \|r_m^{(1)}\|_2 \leq \kappa^{2m} \|r_0\|_2,$$

式中:

$$r_m^{(2)} = b - Ax_m^{(2)}; \quad r_m^{(1)} = b - Ax_m^{(1)} = r_m.$$

如此可证, 若重新开始了 ℓ 次, 产生了 $x_m^{(\ell)}$, 则有

$$\|r_m^{(\ell)}\|_2 \leq \kappa^{\ell m} \|r_0\|_2,$$

式中: $r_m^{(\ell)} = b - Ax_m^{(\ell)}$. 将 κ 的表达式代入上式即得定理的结论. 证毕. \square

注 4.3 定理 4.7 表明, 如果系数矩阵 A 是正定的, 则对任意给定的正数 $m > 0$, GMRES(m) 总是收敛的.

2. 儒科夫斯基映射

要研究 GMRES 方法对更广一类线性方程组的收敛性, 需要借助复变函数中的儒科夫斯基 (Joukowski) 映射来导出复变元的 Chebyshev 多项式的一种易于使用的定义.

儒科夫斯基映射是指如下定义的从复平面到复平面的映射:

$$z = \frac{1}{2}(w + w^{-1}) \equiv J(w). \quad (4.76)$$

对任意的 $w = re^{i\theta}$, $r > 1$, 有

$$z = J(w) = x + iy \equiv a_r \cos \theta + i b_r \sin \theta,$$

式中:

$$a_r = \frac{1}{2}(r + r^{-1}), \quad b_r = \frac{1}{2}(r - r^{-1}). \quad (4.77)$$

从几何上看, J 将 w 平面内的圆

$$C_r = \{w = re^{i\theta} : 0 \leq \theta \leq 2\pi\},$$

映射到 z 平面的一个椭圆

$$E_r = \left\{z = x + iy : \frac{x^2}{a_r^2} + \frac{y^2}{b_r^2} = 1\right\}.$$

显然, 这一映射也将 w 平面内的圆

$$C_{r^{-1}} = \{w = r^{-1}e^{-i\theta} : 0 \leq \theta \leq 2\pi\},$$

映射到了 E_r .

当 $r = 1$ 时, 有 $b_r = 0$. 此时 $z = J(w) = \cos \theta$, 即 J 把单位圆周映射到 z 平面内的线段 $[-1, 1]$. 当 θ 从 0 变到 π 时, z 从 1 变到 -1; 而当 θ 从 π 变到 2π 时, z 从 -1 变到 1.

当 C_r 连续地收缩到单位圆周 C_1 时, E_r 将连续地收缩到线段 $[-1, 1]$. 因此, 由 E_r 所包围的区域

$$E_r = \left\{ z = x + iy : \frac{x^2}{a_r^2} + \frac{y^2}{b_r^2} \leq 1 \right\}$$

中的每一个点 z , 在圆环

$$C_{1,r} = \{ w : 1 \leq |w| \leq r \}$$

中都存在一个点 w , 使得 $z = J(w)$, 即 $J(C_{1,r}) = E_r$.

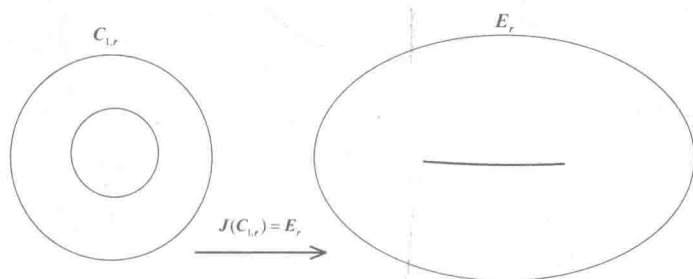


图 4.8 儒可夫斯基映射 $z = J(w)$

现在将 z 平面沿线段 $[-1, 1]$ 切开, 取定 $\sqrt{z^2 - 1}$ 的一个解析分支, 使得

$$w = z + \sqrt{z^2 - 1} \equiv J^{-1}(z). \quad (4.78)$$

刚好在单位圆外, 即 $|w| \geq 1$. 显然, J^{-1} 刚好将 z 平面内的椭圆 E_r 映射到 w 平面内的圆环 $C_{1,r}$, 而且有 $J^{-1}(E_r) = C_{1,r}$.

下面考察 z 平面内的一个中心位于 $z_c = x_c + iy_c$ 的椭圆

$$E(a, b; z_c) = \left\{ z = x + iy : \frac{(x - x_c)^2}{a^2} + \frac{(y - y_c)^2}{b^2} = 1 \right\},$$

这里 $0 < b < a$. 令 $d = \sqrt{a^2 - b^2}$, 即 d 为椭圆的半焦距. 现在作平移伸缩变换

$$\tilde{z} = \frac{z - z_c}{d} \quad \left(\text{即 } \tilde{x} = \frac{x - x_c}{d}, \tilde{y} = \frac{y - y_c}{d} \right),$$

则该变换将椭圆 $E(a, b; z_c)$ 变为中心在原点的椭圆

$$E_{\tilde{r}} = \left\{ \tilde{z} = \tilde{x} + i\tilde{y} : \frac{\tilde{x}^2}{a_{\tilde{r}}^2} + \frac{\tilde{y}^2}{b_{\tilde{r}}^2} = 1 \right\},$$

式中: \tilde{r} 为

$$\frac{a}{d} = \frac{1}{2}(r + r^{-1})$$

的最大根

$$\tilde{r} = \frac{a}{d} + \sqrt{\left(\frac{a}{d}\right)^2 - 1}, \quad (4.79)$$

$a_{\tilde{r}}$ 和 $b_{\tilde{r}}$ 是在式 (4.77) 中将 r 换作 \tilde{r} 而得到的.

再由儒科夫斯基映射的性质, 映射

$$\frac{z - z_c}{d} = \frac{w + w^{-1}}{2} \quad (4.80)$$

实现了区域

$$E(a, b; z_c) = \left\{ z = x + iy : \frac{(x - x_c)^2}{a^2} + \frac{(y - y_c)^2}{b^2} \leq 1 \right\} \quad (4.81)$$

与圆环

$$C_{1, \tilde{r}} = \{ w : 1 \leq |w| \leq \tilde{r} \} \quad (4.82)$$

之间的点与点之间的一一对应关系.

3. 复变元的 Chebyshev 多项式

借助儒科夫斯基映射, 定义函数

$$C_k(z) = \frac{1}{2}(w^k + w^{-k}), \quad (4.83)$$

式中:

$$z = \frac{1}{2}(w + w^{-1}).$$

由定义显然有

$$C_0(z) = \frac{1}{2}(w^0 + w^0) = 1,$$

$$C_1(z) = \frac{1}{2}(w + w^{-1}) = z.$$

而且容易导出

$$\begin{aligned} zC_k(z) &= \frac{1}{2}(w + w^{-1}) \cdot \frac{1}{2}(w^k + w^{-k}) \\ &= \frac{1}{4}[(w^{k+1} + w^{-(k+1)}) + (w^{k-1} + w^{-(k-1)})] \\ &= \frac{1}{2}[C_{k+1}(z) + C_{k-1}(z)], \end{aligned}$$

从而有

$$C_{k+1}(z) = 2zC_k(z) - C_{k-1}(z).$$

这表明由式 (4.83) 所定义的复变函数就是第一类 Chebyshev 多项式. 换句话说, 式 (4.83) 是复变元的 Chebyshev 多项式的又一种表达方式.

再由式 (4.78), 即知式 (4.83) 又可写为

$$C_k(z) = \frac{1}{2} \left[\left(z + \sqrt{z^2 - 1} \right)^k + \left(z + \sqrt{z^2 - 1} \right)^{-k} \right].$$

4. GMRES 方法的收敛性定理

现在考虑 GMRES 方法的收敛性和误差估计. 设 A 是可对角化的, 即存在一个非奇异矩阵 $X \in \mathbb{R}^{n \times n}$, 使得 $A = X \Lambda X^{-1}$, 其中 $\Lambda = \text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_n\}$. 注意到

$$\psi(A) = X\psi(\Lambda)X^{-1}, \quad \psi(\Lambda) = \text{diag}\{\psi(\lambda_1), \psi(\lambda_2), \dots, \psi(\lambda_n)\},$$

由式 (4.67), 得

$$\|r_k\|_2 \leq \|X\|_2 \cdot \|X^{-1}\|_2 \min_{\psi \in \mathcal{P}_k^0} \max_{1 \leq i \leq n} |\psi(\lambda_i)| \cdot \|r_0\|_2. \quad (4.84)$$

再定义

$$\nu(A) = \inf\{\|X\|_2 \|X^{-1}\|_2 : A = X \Lambda X^{-1}\}, \quad (4.85)$$

则有

$$\|r_k\|_2 \leq \nu(A) \min_{\psi \in \mathcal{P}_k^0} \max_{1 \leq i \leq n} |\psi(\lambda_i)| \cdot \|r_0\|_2, \quad (4.86)$$

其中由式 (4.85) 所定义的数 $\nu(A)$ 被称为 A 的谱条件数.

下面通过选择特殊的多项式 $\psi \in \mathcal{P}_k^0$ 给出

$$\min_{\psi \in \mathcal{P}_k^0} \max_{1 \leq i \leq n} |\psi(\lambda_i)|$$

的尽可能小的上界估计. 由于这里的 λ_i 可能是复数, 因此比对称矩阵时的相应估计要困难得多, 需要用到复变元 Chebyshev 多项式的有关性质.

定理 4.8 设式 (4.51) 的系数矩阵 A 是可对角化的, 并满足

$$\lambda(A) \subset E(a, b; z_c),$$

其中 $z_c = c$, 而且 $0 < b < a < c$, 并假设 GMRES 方法已经进行了 k 步得到了 x_k , 则残量 $r_k = b - Ax_k$ 满足

$$\frac{\|r_k\|_2}{\|r_0\|_2} \leq 2\nu(A) \left(\frac{a+b}{c + \sqrt{c^2 + b^2 - a^2}} \right)^k, \quad (4.87)$$

这里的 $\nu(A)$ 由式 (4.85) 所定义.

证明 令 $d = \sqrt{a^2 - b^2}$, 并定义

$$\hat{\psi}(z) = C_k \left(\frac{z - z_c}{d} \right) / C_k \left(\frac{-z_c}{d} \right),$$

则有 $\hat{\psi} \in \mathcal{P}_k^0$. 这样, 由式 (4.86), 有

$$\frac{\|r_k\|_2}{\|r_0\|_2} \leq \nu(A) \max_{\lambda \in \lambda(A)} |\hat{\psi}(\lambda)| \leq \nu(A) \max_{z \in E(a, b, z_c)} |\hat{\psi}(z)|. \quad (4.88)$$

由式 (4.83) 和式 (4.88), 有

$$C_k \left(\frac{z - z_c}{d} \right) = \frac{1}{2} (w^k + w^{-k}). \quad (4.89)$$

式中:

$$\frac{z - z_c}{d} = \frac{1}{2}(w + w^{-1}). \quad (4.90)$$

注意: 式 (4.90) 所定义的映射实现了区域 $E(a, b, z_c)$ 与圆环 $C_{1, \tilde{r}}$ 之间的点与点之间的一一对应关系, 有

$$\max_{z \in E(a, b, z_c)} |\hat{\psi}(z)| = \max_{w \in C_{1, \tilde{r}}} \left| \frac{w^k + w^{-k}}{2C_k(-z_c/d)} \right|. \quad (4.91)$$

在式 (4.90) 中令 $z = 0$, 得

$$-\frac{z_c}{d} = -\frac{c}{d} = \frac{1}{2}(w_0 + w_0^{-1}). \quad (4.92)$$

解此方程, 并选择其模最大者, 即

$$w_0 = -\frac{c}{d} - \sqrt{\left(\frac{c}{d}\right)^2 - 1}. \quad (4.93)$$

此外, 对任意的 $w = re^{i\theta}$, 其中 $1 \leq r \leq \tilde{r}$, 有

$$w^k + w^{-k} = (r^k + r^{-k}) \cos k\theta + i(r^k - r^{-k}) \sin k\theta,$$

从而有

$$\begin{aligned} |w^k + w^{-k}|^2 &= (r^k + r^{-k})^2 \cos^2 k\theta + (r^k - r^{-k})^2 \sin^2 k\theta \\ &= (r^k - r^{-k})^2 + 4 \cos^2 k\theta. \end{aligned}$$

显然, 该函数在圆周 $w = re^{i\theta}$ 上的最大值为

$$(r^k - r^{-k})^2 + 4 = (r^k + r^{-k})^2,$$

从而

$$|w^k + w^{-k}| \leq r^k + r^{-k}, \quad w = re^{i\theta}. \quad (4.94)$$

这样, 由式 (4.89), 式 (4.91), 式 (4.92) 和式 (4.94), 有

$$\max_{z \in E(a, b, z_c)} |\hat{\psi}(z)| = \frac{\tilde{r}^k + \tilde{r}^{-k}}{|w_0^k + w_0^{-k}|}, \quad (4.95)$$

这里 \tilde{r} 由 (4.79) 定义.

此外, 由 \tilde{r} 的定义式 (4.79), w_0 的定义 (4.93) 以及 $d^2 = a^2 - b^2$, 得

$$\tilde{r} = \frac{a+b}{d}, \quad w_0 = -\frac{c + \sqrt{c^2 + b^2 - a^2}}{d},$$

从而

$$\frac{\tilde{r}^k + \tilde{r}^{-k}}{|w_0^k + w_0^{-k}|} = \frac{[(a+b)/d]^k + [(a+b)/d]^{-k}}{[(c + \sqrt{c^2 + b^2 - a^2})/d]^k + [(c + \sqrt{c^2 + b^2 - a^2})/d]^{-k}}$$

$$\begin{aligned}
&= \left(\frac{a+b}{c+\sqrt{c^2+b^2-a^2}} \right)^k \frac{1+[d/(a+b)]^{2k}}{1+[(c+\sqrt{c^2+b^2-a^2})/d]^{-2k}} \\
&\leq \left(\frac{a+b}{c+\sqrt{c^2+b^2-a^2}} \right)^k \left[1 + \left(\frac{a-b}{a+b} \right)^k \right] \\
&\leq 2 \left(\frac{a+b}{c+\sqrt{c^2+b^2-a^2}} \right)^k.
\end{aligned} \tag{4.96}$$

将式 (4.96), 式 (4.95) 和式 (4.88) 结合起来即得定理的结论. 证毕. \square

4.3 极小残量法

本节介绍求解对称不定线性方程组

$$Ax = b \tag{4.97}$$

的极小残量法, 其中系数矩阵 A 是给定的对称不定矩阵 (其特征值有正有负), $b \in \mathbb{R}^n$ 是给定的列向量, 而 $x \in \mathbb{R}^n$ 是待求的未知向量.

若将 4.1 节所介绍的共轭梯度法应用到这类线性方程组上, 其投影方程组 (4.13) 的系数矩阵 T_k 就可能是奇异的, 从而导致算法失败 (即在没有求得式 (4.97) 的很好的近似解之前就发生中断). 基于这种情况, Paige 和 Saunders 于 1975 年在文献 [26] 中提出了一种克服这一困难的方法, 导出了求解对称不定线性方程组的 SYMMLQ 方法. 此外, 在该文献中还给出了求解这类方程组的极小残量法 (MINRES). 本节着重介绍极小残量法, SYMMLQ 方法将在后面再做介绍.

4.3.1 MINRES 方法

极小残量法的出发点是寻找一个 $x_k \in x_0 + \mathcal{K}_k(A, r_0)$, 使得

$$\|r_k\|_2 = \min\{\|b - Ax\|_2 : x \in x_0 + \mathcal{K}_k(A, r_0)\}, \tag{4.98}$$

这里 $r_k = b - Ax_k$ 是残量, 而 $\mathcal{K}_k(A, r_0)$ 是由矩阵 A 和向量 r_0 所生成的 Krylov 子空间.

假设已经求得一个以 $v_1 = r_0/\|r_0\|_2$ 为初始向量的长度为 k 的 Lanczos 分解

$$AV_k = V_{k+1}\tilde{T}_k, \tag{4.99}$$

式中:

$$V_{k+1} = [V_k, v_{k+1}] = [v_1, v_2, \dots, v_k, v_{k+1}] \in \mathbb{R}^{n \times (k+1)}$$

满足 $\mathbf{V}_{k+1}^T \mathbf{V}_{k+1} = \mathbf{I}_{k+1}$, 而 $\tilde{\mathbf{T}}_k$ 是 $(k+1) \times k$ 阶的三对角矩阵, 即

$$\tilde{\mathbf{T}}_k = \begin{bmatrix} \alpha_1 & \beta_1 & & & \\ \beta_1 & \alpha_2 & \beta_2 & & \\ & \ddots & \ddots & \ddots & \\ & & \beta_{k-2} & \alpha_{k-1} & \beta_{k-1} \\ & & & \beta_{k-1} & \alpha_k \\ & & & & 0 & \beta_k \end{bmatrix}, \quad \beta_i \neq 0. \quad (4.100)$$

注意到 $\mathcal{K}_k(\mathbf{A}, \mathbf{r}_0) = \mathcal{R}(\mathbf{V}_k)$, 则对任意的 $\mathbf{x} = \mathbf{x}_0 + \mathbf{V}_k \mathbf{z} \in \mathbf{x}_0 + \mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$, 有

$$\begin{aligned} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2 &= \|\mathbf{r}_0 - \mathbf{A}\mathbf{V}_k \mathbf{z}\|_2 = \|\beta \mathbf{V}_{k+1} \mathbf{e}_1 - \mathbf{V}_{k+1} \tilde{\mathbf{T}}_k \mathbf{z}\|_2 \\ &= \|\mathbf{V}_{k+1}(\beta \mathbf{e}_1 - \tilde{\mathbf{T}}_k \mathbf{z})\|_2 = \|\beta \mathbf{e}_1 - \tilde{\mathbf{T}}_k \mathbf{z}\|_2, \end{aligned} \quad (4.101)$$

式中: $\beta = \|\mathbf{r}_0\|_2$. 这样求 $\mathbf{x}_k \in \mathbf{x}_0 + \mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$ 满足式 (4.98) 就等价于求 $\mathbf{z}_k \in \mathbb{R}^k$, 使得

$$\|\beta \mathbf{e}_1 - \tilde{\mathbf{T}}_k \mathbf{z}_k\|_2 = \min\{\|\beta \mathbf{e}_1 - \tilde{\mathbf{T}}_k \mathbf{z}\|_2 : \mathbf{z} \in \mathbb{R}^k\}. \quad (4.102)$$

一旦这样的 \mathbf{z}_k 求得, 则所寻求的 \mathbf{x}_k 就是 $\mathbf{x}_k = \mathbf{x}_0 + \mathbf{V}_k \mathbf{z}_k$.

极小化问题 (4.102) 是一个系数矩阵为三对角矩阵的最小二乘问题, 现在已有很多十分成熟的数值方法来求它的解. 这里将用 QR 分解来求解式 (4.102).

由于 $\tilde{\mathbf{T}}_k$ 具有式 (4.100) 所示形状, 故可以计算 k 个 Givens 变换 $\mathbf{G}_1, \mathbf{G}_2, \dots, \mathbf{G}_k$, 使得

$$\mathbf{G}_k \mathbf{G}_{k-1} \cdots \mathbf{G}_2 \mathbf{G}_1 \tilde{\mathbf{T}}_k = \begin{bmatrix} \mathbf{R}_k \\ \mathbf{0} \end{bmatrix}, \quad (4.103)$$

式中:

$$\begin{aligned} \mathbf{G}_i &= \text{diag} \left(\mathbf{I}_{i-1}, \begin{bmatrix} c_i & s_i \\ -s_i & c_i \end{bmatrix}, \mathbf{I}_{k-i} \right) \in \mathbb{R}^{(k+1) \times (k+1)}, \quad c_i^2 + s_i^2 = 1, \\ \mathbf{R}_k &= \begin{bmatrix} \gamma_1 & \delta_1 & \varepsilon_1 & & \\ & \gamma_2 & \delta_2 & \ddots & \\ & & \ddots & \ddots & \varepsilon_{k-2} \\ & & & \gamma_{k-1} & \delta_{k-1} \\ & & & & \gamma_k \end{bmatrix}, \end{aligned} \quad (4.104)$$

而且 $\beta_i \neq 0$ 蕴涵着 $\gamma_i \neq 0, i = 1, 2, \dots, k$, 从而 \mathbf{R}_k 是非奇异的. 令

$$\mathbf{G} = \mathbf{G}_k \mathbf{G}_{k-1} \cdots \mathbf{G}_2 \mathbf{G}_1, \quad \begin{bmatrix} \mathbf{t}_k \\ \rho_k \end{bmatrix} = \mathbf{G}(\beta \mathbf{e}_1), \quad \mathbf{t}_k = (\tau_1, \tau_2, \dots, \tau_k)^T, \quad (4.105)$$

则 G 是 $k+1$ 阶正交矩阵, 且 τ_k, ρ_k 如式 (4.58) 所定义, 即

$$\begin{cases} \tau_1 = \beta c_1, \\ \tau_i = (-1)^{i-1} \beta s_1 s_2 \cdots s_{i-1} c_i, \quad i = 2, 3, \cdots, k, \\ \rho_k = (-1)^k \beta s_1 s_2 \cdots s_k. \end{cases} \quad (4.106)$$

这样, 利用分解式 (4.103) 和记号 (4.105), 有

$$\begin{aligned} \|\tilde{T}_k z - \beta e_1\|_2^2 &= \|G(\tilde{T}_k z - \beta e_1)\|_2^2 = \left\| \begin{bmatrix} R_k \\ 0 \end{bmatrix} z - \begin{bmatrix} t_k \\ \rho_k \end{bmatrix} \right\|_2^2 \\ &= \|R_k z - t_k\|_2^2 + \rho_k^2. \end{aligned}$$

对任意的 $z \in \mathbb{R}^k$ 成立. 由此立即知道, 最小二乘问题 (4.102) 有唯一解, 即

$$z_k = R_k^{-1} t_k, \quad (4.107)$$

而且有

$$\|\tilde{T}_k z_k - \beta e_1\|_2 = |\rho_k|. \quad (4.108)$$

由式 (4.107) 求得 z_k 之后, 就可算出所求的 x_k 为 $x_k = x_0 + V_k z_k$.

由式 (4.101) 和式 (4.108), 得

$$\|r_k\|_2 = \|b - Ax_k\|_2 = \|\tilde{T}_k z_k - \beta e_1\|_2 = |\rho_k|,$$

故在实际计算时, 可以用

$$|\rho_k|/\beta \leq \varepsilon \quad (4.109)$$

作为迭代终止的准则, 其中 $\varepsilon > 0$ 是给定的误差要求.

由式 (4.106) 可知, ρ_k 的值并不需要 z_k 和 x_k 的信息. 因此, 只需在 ρ_k 满足式 (4.109) 之后, 再去计算 z_k 和 x_k 即可.

至此, 已经给出了求解极小化问题 (4.98) 的具体方法, 然而这样做的一个最大缺点就是需要保存 V_k 的所有列向量, 随着 k 的增加, 存储量的需求会越来越大. 幸运的是, 在文献 [26] 中还给出了一种巧妙的方法来避免这种情形的出现. 这一方法的具体做法是不将 z_k 明确求出, 而是由 z_k 直接导出计算 x_k 的递推公式.

将式 (4.107) 代入 $x_k = x_0 + V_k z_k$, 得

$$x_k = x_0 + V_k R_k^{-1} t_k = x_0 + P_k t_k, \quad (4.110)$$

式中: $P_k = V_k R_k^{-1}$. 这样, 只要将 P_k 算出, 就可以通过式 (4.110) 计算 x_k . 令 $P_k = [p_1, p_2, \cdots, p_k]$, 比较 $P_k R_k = V_k$ 两边的每一列, 得

$$\begin{aligned} \gamma_1 p_1 &= v_1, \\ \delta_1 p_1 + \gamma_2 p_2 &= v_2, \\ \varepsilon_{i-2} p_{i-2} + \delta_{i-1} p_{i-1} + \gamma_i p_i &= v_i, \quad i = 3, 4, \cdots, k. \end{aligned}$$

由此可求得 P_k 的列向量为

$$\begin{cases} p_1 = v_1/\gamma_1, \\ p_2 = (v_2 - \delta_1 p_1)/\gamma_2, \\ p_i = (v_i - \varepsilon_{i-2} p_{i-2} - \delta_{i-1} p_{i-1})/\gamma_i, \quad i = 3, 4, \dots, k. \end{cases} \quad (4.111)$$

下面借助式 (4.106), 式 (4.110) 和式 (4.111) 导出计算 x_k 的递推公式.

首先注意, 若 Lanczos 分解的长度由 k 增加到 $k+1$, 则有

$$\tilde{T}_{k+1} = \left[\begin{array}{c|c} \tilde{T}_k & \tilde{t}_{k+1} \\ \hline 0 & \beta_{k+1} \end{array} \right], \quad \tilde{t}_{k+1} = (0, \dots, 0, \beta_k, \alpha_{k+1})^T,$$

于是, 有

$$R_{k+1} = \left[\begin{array}{c|c} R_k & \tilde{r}_{k+1} \\ \hline 0 & \gamma_{k+1} \end{array} \right], \quad \tilde{r}_{k+1} = (0, \dots, 0, \varepsilon_{k-1}, \delta_k)^T, \quad (4.112)$$

式中:

$$\begin{aligned} \varepsilon_{k-1} &= s_{k-1} \beta_k, & \tilde{\beta}_k &= c_{k-1} \beta_k, \\ \delta_k &= c_k \tilde{\beta}_k + s_k \alpha_{k+1}, & \tilde{\alpha}_{k+1} &= -s_k \tilde{\beta}_k + c_k \alpha_{k+1}, \end{aligned} \quad (4.113)$$

上面的四个等式由如下 Givens 变换得到:

$$\begin{aligned} \tilde{T}_{k+1} &= \left[\begin{array}{cccccc|c} \alpha_1 & \beta_1 & & & & & \\ & \beta_1 & \alpha_2 & \beta_2 & & & \\ & & \ddots & \ddots & \ddots & & \\ & & & \beta_{k-3} & \alpha_{k-2} & \beta_{k-2} & \\ & & & & \beta_{k-2} & \alpha_{k-1} & \beta_{k-1} \\ & & & & & \beta_{k-1} & \alpha_k & \beta_k \\ & & & & & & \beta_k & \alpha_{k+1} \\ \hline & & & & & & & \beta_{k+1} \end{array} \right] \xrightarrow{G_1} \left[\begin{array}{ccc|c} \gamma_1 & \delta_1 & \varepsilon_1 & \\ & 0 & \tilde{\alpha}_2 & \tilde{\beta}_2 \\ & & \ddots & \ddots \\ & & & \beta_{k-3} & \alpha_{k-2} & \beta_{k-2} \\ & & & & \beta_{k-2} & \alpha_{k-1} & \beta_{k-1} \\ & & & & & \beta_{k-1} & \alpha_k & \beta_k \\ & & & & & & \beta_k & \alpha_{k+1} \\ \hline & & & & & & & \beta_{k+1} \end{array} \right] \\ &\xrightarrow{G_2 \dots G_{k-2}} \left[\begin{array}{ccc|c} \gamma_1 & \delta_1 & \varepsilon_1 & \\ & \gamma_2 & \delta_2 & \varepsilon_2 \\ & & \ddots & \ddots \\ & & & \gamma_{k-2} & \delta_{k-2} & \varepsilon_{k-2} \\ & & & & \tilde{\alpha}_{k-1} & \tilde{\beta}_{k-1} \\ & & & & \beta_{k-1} & \alpha_k & \beta_k \\ & & & & & \beta_k & \alpha_{k+1} \\ \hline & & & & & & \beta_{k+1} \end{array} \right] \end{aligned}$$

$$\xrightarrow{G_{k-1}} \left[\begin{array}{ccc|ccc} \gamma_1 & \delta_1 & \varepsilon_1 & & & \\ & \gamma_2 & \delta_2 & \varepsilon_2 & & \\ & & \ddots & \ddots & \ddots & \\ & & & \gamma_{k-2} & \delta_{k-2} & \varepsilon_{k-2} \\ & & & & \gamma_{k-1} & \delta_{k-1} & \varepsilon_{k-1} \\ & & & & & \tilde{\alpha}_k & \tilde{\beta}_k \\ & & & & & \beta_k & \alpha_{k+1} \\ \hline & & & & & & \beta_{k+1} \end{array} \right] \xrightarrow{G_k} \left[\begin{array}{ccc|ccc} \gamma_1 & \delta_1 & \varepsilon_1 & & & \\ & \gamma_2 & \delta_2 & \varepsilon_2 & & \\ & & \ddots & \ddots & \ddots & \\ & & & \gamma_{k-2} & \delta_{k-2} & \varepsilon_{k-2} \\ & & & & \gamma_{k-1} & \delta_{k-1} & \varepsilon_{k-1} \\ & & & & & \gamma_k & \delta_k \\ & & & & & & \tilde{\alpha}_{k+1} \\ \hline & & & & & & \beta_{k+1} \end{array} \right],$$

即由

$$\begin{bmatrix} c_{k-1} & s_{k-1} \\ -s_{k-1} & c_{k-1} \end{bmatrix} \begin{bmatrix} 0 \\ \beta_k \end{bmatrix} = \begin{bmatrix} \varepsilon_{k-1} \\ \tilde{\beta}_k \end{bmatrix}, \quad \begin{bmatrix} c_k & s_k \\ -s_k & c_k \end{bmatrix} \begin{bmatrix} \tilde{\beta}_k \\ \alpha_{k+1} \end{bmatrix} = \begin{bmatrix} \delta_k \\ \tilde{\alpha}_{k+1} \end{bmatrix}$$

得到. 而 γ_{k+1} 是在确定第 $k+1$ 个 Givens 变换 G_{k+1} 时得到的, 即计算 c_{k+1} 和 s_{k+1} , 使得

$$\begin{bmatrix} c_{k+1} & s_{k+1} \\ -s_{k+1} & c_{k+1} \end{bmatrix} \begin{bmatrix} \tilde{\alpha}_{k+1} \\ \beta_{k+1} \end{bmatrix} = \begin{bmatrix} \gamma_{k+1} \\ 0 \end{bmatrix}. \quad (4.114)$$

由式 (4.106) 可知 $\mathbf{t}_{k+1} = (\mathbf{t}_k^T, \tau_{k+1})^T$, 其中

$$\tau_{k+1} = (-1)^k \beta s_1 s_2 \cdots s_k c_{k+1} = \rho_k c_{k+1}, \quad (4.115)$$

而

$$\rho_{k+1} = (-1)^{k+1} \beta s_1 s_2 \cdots s_k s_{k+1} = -\rho_k s_{k+1}. \quad (4.116)$$

由式 (4.111) 和式 (4.112), 有

$$\mathbf{P}_{k+1} = \mathbf{V}_{k+1} \mathbf{R}_{k+1}^{-1} = [\mathbf{V}_k, \mathbf{v}_{k+1}] \left[\begin{array}{c|c} \mathbf{R}_k^{-1} & \mathbf{s}_{k+1} \\ \hline \mathbf{0} & \gamma_{k+1}^{-1} \end{array} \right] = [\mathbf{P}_k, \mathbf{p}_{k+1}], \quad (4.117)$$

其中

$$\mathbf{p}_{k+1} = (\mathbf{v}_{k+1} - \varepsilon_{k-1} \mathbf{p}_{k-1} - \delta_k \mathbf{p}_k) / \gamma_{k+1}. \quad (4.118)$$

这是由于

$$\begin{aligned} \mathbf{I}_{k+1} &= \mathbf{R}_{k+1} \mathbf{R}_{k+1}^{-1} = \left[\begin{array}{c|c} \mathbf{R}_k & \tilde{\mathbf{r}}_{k+1} \\ \hline \mathbf{0} & \gamma_{k+1} \end{array} \right] \left[\begin{array}{c|c} \mathbf{R}_k^{-1} & \mathbf{s}_{k+1} \\ \hline \mathbf{0} & \gamma_{k+1}^{-1} \end{array} \right] \\ &= \left[\begin{array}{c|c} \mathbf{I}_k & \mathbf{R}_k \mathbf{s}_{k+1} + \gamma_{k+1}^{-1} \tilde{\mathbf{r}}_{k+1} \\ \hline \mathbf{0} & 1 \end{array} \right], \end{aligned}$$

由此可得

$$\mathbf{R}_k \mathbf{s}_{k+1} + \gamma_{k+1}^{-1} \tilde{\mathbf{r}}_{k+1} = \mathbf{0} \implies \mathbf{s}_{k+1} = -\gamma_{k+1}^{-1} \mathbf{R}_k^{-1} \tilde{\mathbf{r}}_{k+1}.$$

于是

$$\begin{aligned} \mathbf{p}_{k+1} &= \mathbf{V}_k \mathbf{s}_{k+1} + \gamma_{k+1}^{-1} \mathbf{v}_{k+1} = -\gamma_{k+1}^{-1} \mathbf{V}_k \mathbf{R}_k^{-1} \tilde{\mathbf{r}}_{k+1} + \gamma_{k+1}^{-1} \mathbf{v}_{k+1} \\ &= (\mathbf{v}_{k+1} - \mathbf{P}_k \tilde{\mathbf{r}}_{k+1}) / \gamma_{k+1} = (\mathbf{v}_{k+1} - \varepsilon_{k-1} \mathbf{p}_{k-1} - \delta_k \mathbf{p}_k) / \gamma_{k+1}. \end{aligned}$$

从而有

$$\mathbf{x}_{k+1} = \mathbf{x}_0 + \mathbf{P}_{k+1} \mathbf{t}_{k+1} = \mathbf{x}_0 + [\mathbf{P}_k, \mathbf{p}_{k+1}] \begin{bmatrix} \mathbf{t}_k \\ \tau_{k+1} \end{bmatrix} = \mathbf{x}_k + \tau_{k+1} \mathbf{p}_{k+1}. \quad (4.119)$$

这就有人们希望得到的 \mathbf{x}_k 的递推公式. 这样, 每次迭代只需保存 $\mathbf{x}_k, \mathbf{p}_{k-1}, \mathbf{p}_k, \mathbf{v}_k, \mathbf{v}_{k+1}$ 这几个向量即可.

综述上面的讨论, 就得到了如下的极小残量方法.

算法 4.9 (MINRES 方法) 给定 n 阶非奇异的实对称矩阵 \mathbf{A} , n 维向量 \mathbf{b} 和允许误差 $\varepsilon > 0$. 本算法计算近似解向量 \mathbf{x}_k , 使得 $\|\mathbf{r}_k\|_2 / \|\mathbf{r}_0\|_2 \leq \varepsilon$, 其中 $\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k$.

选取 \mathbf{x}_0 ; $\mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{x}_0$; $\beta = \|\mathbf{r}_0\|_2$; $\mathbf{v}_1 = \mathbf{r}_0 / \beta$;

$\alpha_1 = \mathbf{v}_1^T \mathbf{A} \mathbf{v}_1$; $\mathbf{u} = \mathbf{A} \mathbf{v}_1 - \alpha_1 \mathbf{v}_1$; $\beta_1 = \|\mathbf{u}\|_2$;

if $\beta_1 = 0$

$\mathbf{x}_1 = \mathbf{x}_0 + \mathbf{r}_0 / \alpha_1$; 结束

else

$\mathbf{v}_2 = \mathbf{u} / \beta_1$;

end

$c_0 = 1$; $s_0 = 0$; $\mathbf{p}_0 = \mathbf{0}$;

确定 $c_1 = \cos \theta_1$ 和 $s_1 = \sin \theta_1$, 使得

$$\begin{bmatrix} c_1 & s_1 \\ -s_1 & c_1 \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \beta_1 \end{bmatrix} = \begin{bmatrix} \gamma_1 \\ 0 \end{bmatrix};$$

$\mathbf{p}_1 = \mathbf{v}_1 / \gamma_1$; $\rho_1 = -\beta s_1$; $\tau_1 = \beta c_1$;

$\mathbf{x}_1 = \mathbf{x}_0 + \tau_1 \mathbf{p}_1$; $k = 1$;

while ($|\rho_k| / \beta > \varepsilon$)

$\alpha_{k+1} = \mathbf{v}_{k+1}^T \mathbf{A} \mathbf{v}_{k+1}$;

$\mathbf{u} = \mathbf{A} \mathbf{v}_{k+1} - \alpha_{k+1} \mathbf{v}_{k+1} - \beta_k \mathbf{v}_k$;

$\beta_{k+1} = \|\mathbf{u}\|_2$;

if $\beta_{k+1} \neq 0$

$\mathbf{v}_{k+2} = \mathbf{u} / \beta_{k+1}$;

end

$\varepsilon_{k-1} = s_{k-1} \beta_k$; $\tilde{\beta}_k = c_{k-1} \beta_k$;

$\delta_k = c_k \tilde{\beta}_k + s_k \alpha_{k+1}$; $\tilde{\alpha}_{k+1} = -s_k \tilde{\beta}_k + c_k \alpha_{k+1}$;

确定 $c_{k+1} = \cos \theta_{k+1}$ 和 $s_{k+1} = \sin \theta_{k+1}$, 使得

$$\begin{bmatrix} c_{k+1} & s_{k+1} \\ -s_{k+1} & c_{k+1} \end{bmatrix} \begin{bmatrix} \tilde{\alpha}_{k+1} \\ \beta_{k+1} \end{bmatrix} = \begin{bmatrix} \gamma_{k+1} \\ 0 \end{bmatrix};$$

$$\begin{aligned}\tau_{k+1} &= \rho_k c_{k+1}; \quad \rho_{k+1} = -\rho_k s_{k+1}; \\ \mathbf{p}_{k+1} &= (\mathbf{v}_{k+1} - \varepsilon_{k-1} \mathbf{p}_{k-1} - \delta_k \mathbf{p}_k) / \gamma_{k+1}; \\ \mathbf{x}_{k+1} &= \mathbf{x}_k + \tau_{k+1} \mathbf{p}_{k+1}; \\ k &= k + 1;\end{aligned}$$

end

MINRES 方法的 MATLAB 程序如下:

```
%MINRES方法的程序-mminres.m
function [x,iter,time,res,resvec]=mminres(A,b,x,max_it,tol)
%极小残量法求解对称不定方程组Ax=b
tic; r=b-A*x; b=norm(r); v=r/b; bt=b;
z=A*v; a=v'*z; v=z-a*v; b1=norm(v);
if (b1~=0), v1=v/b1; end
c0=1; s0=0; p0=zeros(length(b),1);
[c,s,gama]=givens(a,b1); %Givens变换
p=v/gama; rho=-b*s; tau=b*c;
x=x+tau*p; iter=1;
while (iter<max_it)
    res=abs(rho)/bt; resvec(iter)=res;
    if (res<tol), break; end
    z=A*v1; a=v1'*z; v=z-a*v1-b1*v; b2=norm(v);
    if (b2~=0), v2=v/b2; end
    epsi=s0*b1; bh1=c0*b1;
    dta=c*bh1+s*a; ah=-s*bh1+c*a;
    [c1,s1,gama]=givens(ah,b2); %Givens变换
    tau=rho*c1; rho=-rho*s1;
    p1=(v1-epsi*p0-dta*p)/gama;
    x=x+tau*p1; b1=b2; iter=iter+1;
    v=v1; v1=v2; p0=p; p=p1;
    c0=c; c=c1; s0=s; s=s1;
end
time=toc;
```

例 4.8 假设线性方程组 $\mathbf{Ax} = \mathbf{b}$, 其中矩阵 \mathbf{A} 来自

<http://www.cise.ufl.edu/research/sparse/matrices/PARSECS/SiNa.html>,

它是一个 5743 阶的实对称不定稀疏矩阵, 右端项 \mathbf{b} 为

$$\mathbf{b} = \mathbf{A} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \in \mathbb{R}^n.$$

显然, 该方程组的真解为 $\mathbf{x}^* = (1, 1, \dots, 1)^T$. 应用算法 4.9 到该线性方程组上, 第 224 步得到的 $\hat{\mathbf{x}} = \mathbf{x}_{224}$ 满足

$$\|\hat{\mathbf{x}} - \mathbf{x}^*\|_2 = 5.9139 \times 10^{-8},$$

迭代过程的收敛轨迹如图 4.9 所示, 其中横坐标为迭代步数 k , 纵坐标为相对残差 $\|\mathbf{r}_k\|_2/\|\mathbf{r}_0\|_2$, 这里 \mathbf{r}_k 是第 k 步得到的残差向量.

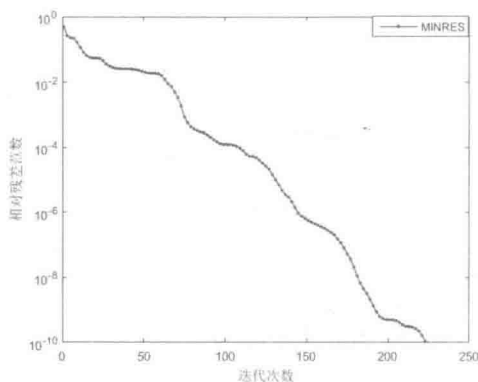


图 4.9 MINRES 方法的收敛特性

4.3.2 PMINRES 方法

MINRES 方法 (算法 4.9) 与 4.1 节的共轭梯度法在形式上很不相同, 且不利于导出相应的预处理算法. 下面从另一个途径推导出与共轭梯度法形式类似的 MINRES 方法, 并由此导出预处理极小残量法 (简记为 PMINRES 方法) 的递推算法. 由前面的论述可知, 极小残量法就是求向量 \mathbf{z}_k 及 $\mathbf{x}_k = \mathbf{x}_0 + \mathbf{V}_k \mathbf{z}_k$ 使得 $\|\mathbf{r}_k\|_2^2$ 达到极小, 其中 $\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k$. 由于 $\mathcal{K}_k(\mathbf{A}, \mathbf{r}_0) = \text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$, 因此, 若记 $\mathbf{S}_k = [\mathbf{r}_0, \mathbf{A}\mathbf{r}_0, \dots, \mathbf{A}^{k-1}\mathbf{r}_0]$, 则 \mathbf{x}_k 也可表示成

$$\mathbf{x}_k = \mathbf{x}_0 + \mathbf{S}_k \mathbf{z}_k.$$

于是问题等价于求 \mathbf{z}_k 使 $\|\mathbf{r}_k\|_2^2$ 达到极小. 注意到

$$\begin{aligned} \|\mathbf{r}_k\|_2^2 &= (\mathbf{A}(\mathbf{x}^* - \mathbf{x}_k), \mathbf{A}(\mathbf{x}^* - \mathbf{x}_k)) \\ &= (\mathbf{A}^T \mathbf{A}(\mathbf{x}^* - \mathbf{x}_0 - \mathbf{S}_k \mathbf{z}_k), \mathbf{x}^* - \mathbf{x}_0 - \mathbf{S}_k \mathbf{z}_k). \end{aligned} \quad (4.120)$$

因为 $\mathbf{A}^T \mathbf{A} = \mathbf{A}^2$ 是对称正定矩阵, 因此可以定义一种新内积

$$(\mathbf{x}, \mathbf{z})_{\mathbf{A}^2} = (\mathbf{A}^T \mathbf{A} \mathbf{x}, \mathbf{z}),$$

因此

$$\|\mathbf{r}_k\|_2^2 = (\mathbf{x}^* - \mathbf{x}_0 - \mathbf{S}_k \mathbf{z}_k, \mathbf{x}^* - \mathbf{x}_0 - \mathbf{S}_k \mathbf{z}_k)_{\mathbf{A}^2}.$$

按照极小与正交关系定理, 可得 \mathbf{z}_k 当且仅当满足条件 (即 $\mathbf{r}_k \perp \mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$)

$$(\mathbf{x}^* - \mathbf{x}_0 - \mathbf{S}_k \mathbf{z}_k, \mathbf{A}^i \mathbf{r}_0)_{\mathbf{A}^2} = 0, \quad i = 0, 1, \dots, k-1.$$

若将 Krylov 子空间 $\text{span}\{\mathbf{r}_0, A\mathbf{r}_0, \dots, A^{k-1}\mathbf{r}_0\}$ 按内积 $(\cdot, \cdot)_{A^2}$ 标准正交化得 $\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_{k-1}$, 则有

$$\mathbf{P}_k = [\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_{k-1}] = \mathbf{S}_k \mathbf{R}_k,$$

式中: \mathbf{R}_k 为上三角矩阵, 它在 $\mathbf{r}_0, A\mathbf{r}_0, \dots, A^{k-1}\mathbf{r}_0$ 线性无关的假设下是非奇异的, 且在对角元确定后是唯一确定的. 于是有

$$\mathbf{S}_k \mathbf{z}_k = \mathbf{P}_k \mathbf{R}_k^{-1} \mathbf{z}_k.$$

记 $\alpha_k := \mathbf{R}_k^{-1} \mathbf{z}_k$, 则

$$\mathbf{x}_k = \mathbf{x}_0 + \mathbf{P}_k \alpha_k.$$

注意到

$$\mathbf{P}_k \alpha_k = \sum_{j=0}^{k-1} \alpha_j^{(k)} \mathbf{p}_j.$$

由

$$(\mathbf{x}^* - \mathbf{x}_0 - \mathbf{P}_k \alpha_k, \mathbf{p}_i)_{A^2} = 0, \quad i = 0, 1, \dots, k-1$$

及 \mathbf{p}_i 的 $(\cdot, \cdot)_{A^2}$ 正交性, 得

$$\begin{aligned} \alpha_i^{(k)} &= \frac{(\mathbf{x}^* - \mathbf{x}_0, \mathbf{p}_i)_{A^2}}{(\mathbf{p}_i, \mathbf{p}_i)_{A^2}} = \frac{(A\mathbf{x}^* - A\mathbf{x}_0, A\mathbf{p}_i)}{(A\mathbf{p}_i, A\mathbf{p}_i)} \\ &= \frac{(b - A\mathbf{x}_0, A\mathbf{p}_i)}{(A\mathbf{p}_i, A\mathbf{p}_i)} = \frac{(\mathbf{r}_0, A\mathbf{p}_i)}{(A\mathbf{p}_i, A\mathbf{p}_i)}. \end{aligned}$$

注意到 $\alpha_i^{(k)}$ 计算公式与 k 无关, 于是可以得到简单的递推公式, 即

$$\begin{aligned} \mathbf{x}_k &= \mathbf{x}_0 + \mathbf{P}_k \alpha_k = \mathbf{x}_0 + \sum_{i=0}^{k-1} \alpha_i \mathbf{p}_i \\ &= \mathbf{x}_0 + \sum_{i=0}^{k-2} \alpha_i \mathbf{p}_i + \alpha_{k-1} \mathbf{p}_{k-1} \\ &= \mathbf{x}_{k-1} + \alpha_{k-1} \mathbf{p}_{k-1}. \end{aligned}$$

现在剩下的问题是如何求得标准正交向量组 $\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_{k-1}$. 按照共轭梯度法的思想, 取

$$\mathbf{p}_0 = \mathbf{r}_0, \quad \mathbf{p}_1 = A\mathbf{p}_0 + \gamma_0^{(1)} \mathbf{p}_0.$$

要使 $(\mathbf{p}_0, \mathbf{p}_1)_{A^2} = 0$, 必须

$$\gamma_0^{(1)} = -\frac{(A\mathbf{p}_0, \mathbf{p}_0)_{A^2}}{(\mathbf{p}_0, \mathbf{p}_0)_{A^2}} = -\frac{(A^2 \mathbf{p}_0, A\mathbf{p}_0)}{(A\mathbf{p}_0, A\mathbf{p}_0)}.$$

然后选取

$$\mathbf{p}_2 = A\mathbf{p}_1 + \gamma_1^{(2)} \mathbf{p}_1 + \gamma_0^{(2)} \mathbf{p}_0.$$

利用 p_0, p_1, p_2 的正交性, 分别用 p_1, p_0 与上式两边作 $(\cdot, \cdot)_{A^2}$ 内积, 得

$$\begin{aligned}\gamma_1^{(2)} &= -\frac{(Ap_1, p_1)_{A^2}}{(p_1, p_1)_{A^2}} = -\frac{(A^2 p_1, Ap_1)}{(Ap_1, Ap_1)}, \\ \gamma_0^{(2)} &= -\frac{(Ap_1, p_0)_{A^2}}{(p_0, p_0)_{A^2}} = -\frac{(A^2 p_1, Ap_0)}{(Ap_0, Ap_0)}.\end{aligned}$$

同理, 令

$$p_3 = Ap_2 + \gamma_2^{(3)} p_2 + \gamma_1^{(3)} p_1 + \gamma_0^{(3)} p_0.$$

分别用 p_2, p_1, p_0 与上式两边作 $(\cdot, \cdot)_{A^2}$ 内积得

$$\begin{aligned}\gamma_2^{(3)} &= -\frac{(Ap_2, p_2)_{A^2}}{(p_2, p_2)_{A^2}} = -\frac{(A^2 p_2, Ap_2)}{(Ap_2, Ap_2)}, \\ \gamma_1^{(3)} &= -\frac{(Ap_2, p_1)_{A^2}}{(p_1, p_1)_{A^2}} = -\frac{(A^2 p_2, Ap_1)}{(Ap_1, Ap_1)}, \\ \gamma_0^{(3)} &= -\frac{(Ap_2, p_0)_{A^2}}{(p_0, p_0)_{A^2}} = -\frac{(p_2, Ap_0)_{A^2}}{(p_0, p_0)_{A^2}} \\ &= -\frac{(p_2, p_1 - \gamma_0^{(1)} p_0)_{A^2}}{(p_0, p_0)_{A^2}} = 0.\end{aligned}$$

依此类推, 有

$$\begin{aligned}p_{i+1} &= Ap_i + \lambda_i p_i + \mu_i p_{i-1}, \\ \lambda_i &= -\frac{(A^2 p_i, Ap_i)}{(Ap_i, Ap_i)}, \quad \mu_i = -\frac{(A^2 p_i, Ap_{i-1})}{(Ap_{i-1}, Ap_{i-1})}, \\ \mu_0 &= 0, \quad i = 0, 1, 2, \dots\end{aligned}$$

综合上述, 即可得到形式与共轭梯度法类似的极小残量递推算法:

取 $p_0 = r_0, \mu_0 = 0$, 对 $k = 0, 1, 2, \dots$ 进行下列计算直到满足终止条件:

$$\alpha_k = \frac{(r_0, Ap_k)}{(Ap_k, Ap_k)}, \quad x_{k+1} = x_k + \alpha_k p_k; \quad (4.121)$$

$$\begin{aligned}\lambda_k &= -\frac{(A^2 p_k, Ap_k)}{(Ap_k, Ap_k)}, \quad \mu_k = -\frac{(A^2 p_k, Ap_{k-1})}{(Ap_{k-1}, Ap_{k-1})}, \\ p_{k+1} &= Ap_k + \lambda_k p_k + \mu_k p_{k-1}.\end{aligned} \quad (4.122)$$

上述算法每次迭代需要作两次矩阵与向量乘法 $(Ap_k, A(Ap_k))$, 要存储四个向量 x_k, p_k, p_{k-1}, Ap_k , 计算量和存储量均为共轭梯度法的两倍.

回顾共轭梯度法中, 在求 p_{k+1} 时利用了 r_{k+1} 代替 Ap_k 来减少计算量和存储量. 在此, 也来分析一下 r_{k+1} . 由

$$(x^* - x_0 - P_k \alpha_k, p_i)_{A^2} = 0, \quad i = 0, 1, \dots, k-1,$$

可知

$$(x^* - x_k, p_i)_{A^2} = 0, \quad i = 0, 1, \dots, k-1.$$

上式即

$$(r_k, Ap_i) = (A(x^* - x_k), Ap_i) = 0, \quad i = 0, 1, \dots, k-1. \quad (4.123)$$

但

$$r_{k+1} = r_k - \alpha_k Ap_k, \quad (4.124)$$

即

$$r_k = r_{k+1} + \alpha_k Ap_k,$$

故有

$$\|r_k\|_2^2 = \|r_{k+1}\|_2^2 + \alpha_k^2 \|Ap_k\|_2^2.$$

因此, 只要 $\alpha_k \|Ap_k\|_2 \neq 0$, 就必有

$$\|r_{k+1}\|_2 < \|r_k\|_2.$$

由式 (4.122), 有

$$\begin{aligned} p_i &= Ap_{i-1} + \lambda_{i-1} p_{i-1} + \mu_{i-1} p_{i-2} = \dots \\ &= A^i p_0 + \sum_{s=0}^{i-1} \delta_s^{(i)} A^s p_0 = A^i r_0 + \sum_{s=0}^{i-1} \delta_s^{(i)} A^s r_0, \end{aligned}$$

得

$$\begin{aligned} r_{i+1} &= r_i - \alpha_i Ap_i = r_i - \alpha_i A \left(A^i r_0 + \sum_{s=0}^{i-1} \delta_s^{(i)} A^s r_0 \right) \\ &= r_i - \alpha_i \sum_{s=0}^{i-1} \delta_s^{(i)} A^{s+1} r_0 - \alpha_i A^{i+1} r_0 \\ &= \dots = \mathcal{P}_i(A) r_0 - \alpha_i A^{i+1} r_0, \end{aligned}$$

式中: $\mathcal{P}_i(A)$ 为关于矩阵 A 的 i 次多项式. 由上式, 得

$$[r_0, r_1, \dots, r_{k-1}] = [r_0, Ar_0, \dots, A^{k-1}r_0] \tilde{R}_k = S_k \tilde{R}_k, \quad (4.125)$$

式中: \tilde{R}_k 为 k 阶上三角矩阵, 且其对角元为 $1, -\alpha_0, -\alpha_1, \dots, -\alpha_{k-2}$.

若构造关于内积 $(\cdot, \cdot)_{A^2}$ 的正交化向量组 $\tilde{p}_0, \tilde{p}_1, \dots, \tilde{p}_{k-1}$:

$$\begin{aligned} \tilde{p}_0 &= r_0, \\ \tilde{p}_1 &= r_1 + \gamma_0^{(1)} \tilde{p}_0, \\ \tilde{p}_2 &= r_2 + \gamma_1^{(2)} \tilde{p}_1 + \gamma_0^{(2)} \tilde{p}_0, \\ &\vdots \end{aligned}$$

$$\tilde{\mathbf{p}}_i = \mathbf{r}_i + \sum_{s=0}^{i-1} \gamma_s^{(i)} \tilde{\mathbf{p}}_s = \mathbf{r}_i + \gamma_{i-1}^{(i)} \tilde{\mathbf{p}}_{i-1} + \cdots + \gamma_1^{(i)} \tilde{\mathbf{p}}_1 + \gamma_0^{(i)} \tilde{\mathbf{p}}_0,$$

式中:

$$\gamma_s^{(i)} = -\frac{(\mathbf{r}_i, \tilde{\mathbf{p}}_s)_{A^2}}{(\tilde{\mathbf{p}}_s, \tilde{\mathbf{p}}_s)_{A^2}}, \quad i = 0, 1, \dots, k-1; \quad s = 0, 1, \dots, i-1.$$

不失一般性, 设 $\tilde{\mathbf{p}}_i = \omega_i \mathbf{p}_i$, 利用 $\mathbf{A}^T = \mathbf{A}$, 有

$$\begin{aligned} (\mathbf{r}_i, \tilde{\mathbf{p}}_s)_{A^2} &= (\mathbf{A}\mathbf{r}_i, \mathbf{A}\tilde{\mathbf{p}}_s) = (\mathbf{A}\mathbf{r}_i, \omega_s \mathbf{A}\mathbf{p}_s) \\ &= (\mathbf{A}\mathbf{r}_i, \omega_s (\mathbf{p}_{s+1} - \lambda_s \mathbf{p}_s - \mu_s \mathbf{p}_{s-1})) \\ &= (\mathbf{r}_i, \omega_s \mathbf{A}\mathbf{p}_{s+1} - \omega_s \lambda_s \mathbf{A}\mathbf{p}_s - \omega_s \mu_s \mathbf{A}\mathbf{p}_{s-1}). \end{aligned}$$

由上式及式 (4.123), 不难发现, 当 $s+1 \leq i-1$, 即 $s \leq i-2$ 时, 有

$$(\mathbf{r}_i, \tilde{\mathbf{p}}_s)_{A^2} = 0,$$

即

$$\gamma_s^{(i)} = 0, \quad s \leq i-2, \quad i = 0, 1, \dots, k-1.$$

于是有

$$\tilde{\mathbf{p}}_i = \mathbf{r}_i + \gamma_{i-1}^{(i)} \tilde{\mathbf{p}}_{i-1}, \quad i = 0, 1, \dots, k-1, \quad \tilde{\mathbf{p}}_{-1} = \mathbf{0}.$$

回顾式 (4.122), \mathbf{p}_i 与 $\mathbf{A}\mathbf{p}_{i-1}$, \mathbf{p}_{i-1} , \mathbf{p}_{i-2} 有关, 可以看到, $\tilde{\mathbf{p}}_i$ 只与 \mathbf{r}_i 和 $\tilde{\mathbf{p}}_{i-1}$ 有关. 现记 $\beta_{i-1} := \gamma_{i-1}^{(i)}$, 在不混淆的情况下, 再记 $\mathbf{p}_i := \tilde{\mathbf{p}}_i$, 则有

$$\begin{aligned} \mathbf{p}_i &= \mathbf{r}_i + \beta_{i-1} \mathbf{p}_{i-1}, \quad i = 1, 2, \dots, k-1, \\ \beta_{i-1} &= \frac{(\mathbf{r}_i, \mathbf{p}_{i-1})_{A^2}}{(\mathbf{p}_{i-1}, \mathbf{p}_{i-1})_{A^2}}. \end{aligned}$$

注意到, 对这一组 $\tilde{\mathbf{p}}_0, \tilde{\mathbf{p}}_2, \dots, \tilde{\mathbf{p}}_{k-1}$ (已记为 $\mathbf{p}_0, \mathbf{p}_2, \dots, \mathbf{p}_{k-1}$), 也有

$$[\mathbf{p}_0, \mathbf{p}_2, \dots, \mathbf{p}_{k-1}] = [\mathbf{r}_0, \mathbf{A}\mathbf{r}_0, \dots, \mathbf{A}^{k-1}\mathbf{r}_0] \hat{\mathbf{R}}_k = \mathbf{S}_k \hat{\mathbf{R}}_k, \quad (4.126)$$

式中: $\hat{\mathbf{R}}_k$ 为上三角矩阵, 对角元也为 $1, -\alpha_0, -\alpha_1, \dots, -\alpha_{k-2}$. 由式 (4.125) 和式 (4.126), 得

$$[\mathbf{r}_0, \mathbf{r}_2, \dots, \mathbf{r}_{k-1}] = [\mathbf{p}_0, \mathbf{p}_2, \dots, \mathbf{p}_{k-1}] \hat{\mathbf{R}}_k^{-1} \tilde{\mathbf{R}}_k,$$

式中: $\hat{\mathbf{R}}_k^{-1} \tilde{\mathbf{R}}_k$ 仍为上三角矩阵, 且对角元均为 1. 于是有

$$\mathbf{r}_i = \mathbf{p}_i + \sum_{s=0}^{i-1} c_s \mathbf{p}_s,$$

及

$$\beta_{i-1} = \frac{(\mathbf{r}_i, \mathbf{p}_{i-1})_{A^2}}{(\mathbf{p}_{i-1}, \mathbf{p}_{i-1})_{A^2}} = \frac{(\mathbf{A}\mathbf{r}_i, \mathbf{A}\mathbf{p}_{i-1})}{(\mathbf{A}\mathbf{p}_{i-1}, \mathbf{A}\mathbf{p}_{i-1})} = \frac{(\mathbf{A}\mathbf{r}_i, \frac{\mathbf{r}_i - \mathbf{r}_{i-1}}{\alpha_{i-1}})}{(\mathbf{A}\mathbf{p}_{i-1}, \mathbf{A}\mathbf{p}_{i-1})}.$$

由于

$$(Ar_i, r_{i-1}) = \left(r_i, A \left(p_{i-1} + \sum_{s=0}^{i-2} c_s p_s \right) \right) = 0,$$

因此

$$\beta_{i-1} = \frac{(Ar_i, r_i)}{\alpha_{i-1}(Ap_{i-1}, Ap_{i-1})}. \quad (4.127)$$

另外, 由 $x_i = x_{i-1} + \alpha_{i-1}p_{i-1}$ 及式 (4.123), 得

$$\begin{aligned} \alpha_{i-1} &= \frac{(x_i - x_{i-1}, p_{i-1})_{A^2}}{(p_{i-1}, p_{i-1})_{A^2}} = \frac{(A(x_i - x_{i-1}), Ap_{i-1})}{(Ap_{i-1}, Ap_{i-1})} \\ &= \frac{(r_{i-1} - r_i, Ap_{i-1})}{(Ap_{i-1}, Ap_{i-1})} = \frac{(r_{i-1}, Ap_{i-1})}{(Ap_{i-1}, Ap_{i-1})} \\ &= \frac{(Ar_{i-1}, p_{i-1})}{(Ap_{i-1}, Ap_{i-1})} = \frac{(Ar_{i-1}, r_{i-1} + \beta_{i-2}p_{i-2})}{(Ap_{i-1}, Ap_{i-1})} \\ &= \frac{(Ar_{i-1}, r_{i-1})}{(Ap_{i-1}, Ap_{i-1})}. \end{aligned} \quad (4.128)$$

将式 (4.128) 代入式 (4.127), 得

$$\beta_{i-1} = \frac{(Ar_i, r_i)}{(Ar_{i-1}, r_{i-1})}.$$

于是可以得到极小残量法另一种形式的递推算法:

算法 4.10 (递推形式的 MINRES 方法) 给定 n 阶非奇异的实对称矩阵 A , n 维向量 b 和允许误差 $\varepsilon > 0$. 本算法计算近似解向量 x_k , 使得 $\|r_k\|_2 / \|r_0\|_2 \leq \varepsilon$, 其中 $r_k = b - Ax_k$.

选取 x_0 ; 计算 $r_0 = b - Ax_0$; $p_0 = r_0$;

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

$$\alpha_k = \frac{(r_k, Ap_k)}{(Ap_k, Ap_k)}; \quad x_{k+1} = x_k + \alpha_k p_k;$$

$$r_{k+1} = r_k - \alpha_k Ap_k; \quad \beta_k = \frac{(Ar_{k+1}, r_{k+1})}{(Ar_k, r_k)};$$

$$p_{k+1} = r_{k+1} + \beta_k p_k; \quad Ap_{k+1} = Ar_{k+1} + \beta_k Ap_k;$$

end

从形式上看, 算法 4.10 每一步迭代也需要计算两次矩阵与向量的乘法 Ap_k 和 Ar_{k+1} , 实际上可从 $Ap_{k+1} = Ar_{k+1} + \beta_k Ap_k$ 来计算 Ap_{k+1} , 这样就只需要计算一次矩阵与向量的乘法. 此外, 在计算中需要保存五个向量 x_k, p_k, r_k, Ap_k 和 Ar_k . 递推形式的 MINRES 方法的 MATLAB 程序如下:

```

%递推形式的极小残量法-minres1.m
function [x,iter,time,res,resvec]=minres1(A,b,x,max_it,tol)
%极小残量法求解对称不定方程组Ax=b
tic; r=b-A*x; p=r; mr=norm(r);
u=A*p; z=A*r; iter=1;
while (iter<max_it)
    alpha=(r'*u)/(u'*u);
    x=x+alpha*p; r1=r-alpha*u;
    z1=A*r1; beta=(z1'*r1)/(z1'*r1);
    p=r1+beta*p; u=z1+beta*u;
    res=norm(r1)/mr; resvec(iter)=res;
    if (res<tol), break; end
    r=r1; z=z1; iter=iter+1;
end
time=toc;

```

现在考虑预处理的极小残量法. 选择一个适当的对称正定矩阵 M , 设 M 有分解 $M = LL^T$, 则方程组 $Ax = b$ 经过预处理后变为

$$\tilde{A}\tilde{x} = \tilde{b}, \quad (4.129)$$

式中:

$$\tilde{A} = L^{-1}AL^{-T}, \quad \tilde{x} = L^Tx, \quad \tilde{b} = L^{-1}b.$$

方程组 (4.129) 的系数矩阵 $\tilde{A} = L^{-1}AL^{-T}$ 仍为对称不定矩阵. 通过 M 和 L 的选取, 可以使得 $L^{-1}AL^{-T}$ 比 A 有更好的条件数. 现在对方程组 (4.129) 使用极小残量法, 得

$$\begin{aligned} \tilde{r}_0 &= \tilde{b} - \tilde{A}\tilde{x}_0, \quad \tilde{p}_0 = \tilde{r}_0, \\ \alpha_k &= \frac{(\tilde{r}_k, \tilde{A}\tilde{p}_k)}{(\tilde{A}\tilde{p}_k, \tilde{A}\tilde{p}_k)}, \quad \tilde{x}_{k+1} = \tilde{x}_k + \alpha_k\tilde{p}_k, \\ \tilde{r}_{k+1} &= \tilde{r}_k - \alpha_k\tilde{A}\tilde{p}_k, \quad \beta_k = \frac{(\tilde{A}\tilde{r}_{k+1}, \tilde{r}_{k+1})}{(\tilde{A}\tilde{r}_k, \tilde{r}_k)}, \\ \tilde{p}_{k+1} &= \tilde{r}_{k+1} + \beta_k\tilde{p}_k, \quad k = 0, 1, \dots \end{aligned}$$

为了在公式中使用原来变量, 作变换: $x_k = L^{-T}\tilde{x}_k$, $p_k = L^{-T}\tilde{p}_k$, 则 $r_k = b - Ax_k = L\tilde{r}_k$.

$$\begin{aligned} \alpha_k &= \frac{(L^{-1}r_k, (L^{-1}AL^{-T})L^Tp_k)}{(L^{-1}AL^{-T}\tilde{p}_k, L^{-1}AL^{-T}\tilde{p}_k)} = \frac{(L^{-T}L^{-1}r_k, Ap_k)}{(L^{-T}L^{-1}Ap_k, Ap_k)} = \frac{(M^{-1}r_k, Ap_k)}{(M^{-1}Ap_k, Ap_k)}, \\ \tilde{x}_{k+1} = \tilde{x}_k + \alpha_k\tilde{p}_k &\Rightarrow L^Tx_{k+1} = L^Tx_k + \alpha_kL^Tp_k \Rightarrow x_{k+1} = x_k + \alpha_kp_k, \end{aligned}$$

$$L^{-1}r_{k+1} = L^{-1}r_k - \alpha_k(L^{-1}AL^{-T})L^T p_k \Rightarrow r_{k+1} = r_k - \alpha_k A p_k,$$

$$\begin{aligned}\beta_k &= \frac{((L^{-1}AL^{-T})\tilde{r}_{k+1}, \tilde{r}_{k+1})}{((L^{-1}AL^{-T})\tilde{r}_k, \tilde{r}_k)} = \frac{(AL^{-T}\tilde{r}_{k+1}, L^{-T}\tilde{r}_{k+1})}{(AL^{-T}\tilde{r}_k, L^{-T}\tilde{r}_k)} \\ &= \frac{(AL^{-T}L^{-1}r_{k+1}, L^{-T}L^{-1}r_{k+1})}{(AL^{-T}L^{-1}r_k, L^{-T}L^{-1}r_k)} = \frac{(A(M^{-1}r_{k+1}), M^{-1}r_{k+1})}{(A(M^{-1}r_k), M^{-1}r_k)},\end{aligned}$$

$$\begin{aligned}\tilde{p}_{k+1} &= \tilde{r}_{k+1} + \beta_k \tilde{p}_k \Rightarrow L^T p_{k+1} = L^{-1}r_{k+1} + \beta_k L^T p_k \\ &\Rightarrow p_{k+1} = L^{-T}L^{-1}r_{k+1} + \beta_k p_k \Rightarrow p_{k+1} = M^{-1}r_{k+1} + \beta_k p_k.\end{aligned}$$

这样, 就得到了预处理极小残量法的递推算法, 具体步骤为:

算法 4.11 (PMINRES 方法) 给定 n 阶非奇异的实对称矩阵 A , n 维向量 b , 允许误差 $\varepsilon > 0$ 和预处理矩阵 M (对称正定). 本算法计算近似解向量 x_k , 使得 $\|r_k\|_2/\|r_0\|_2 \leq \varepsilon$, 其中 $r_k = b - Ax_k$.

选取 $x_0 \in \mathbb{R}^n$; $r_0 = b - Ax_0$; $z_0 = M^{-1}r_0$;

$p_0 = z_0$; $u_0 = Ap_0$; $w_0 = M^{-1}u_0$; $k = 0$;

while ($\|r_k\|_2/\|r_0\|_2 > \varepsilon$)

$$\alpha_k = \frac{(z_k, u_k)}{(w_k, u_k)}; \quad x_{k+1} = x_k + \alpha_k p_k;$$

$$r_{k+1} = r_k - \alpha_k u_k; \quad z_{k+1} = M^{-1}r_k;$$

$$\beta_k = \frac{(Az_{k+1}, z_{k+1})}{(Az_k, z_k)}; \quad p_{k+1} = z_{k+1} + \beta_k p_k;$$

$$u_{k+1} = Ap_{k+1} = Az_{k+1} + \beta_k u_k; \quad w_{k+1} = M^{-1}u_{k+1};$$

$$k = k + 1;$$

end

注 4.4 算法 4.11 仅与预处理子 M 有关, 而与 L 无关. 当 M 为单位阵时, 它就是没有经过预处理的极小残量法. 与共轭梯度法相比较, 预处理极小残量法增加了计算 z_{k+1} 和 w_{k+1} 的工作量, 这可以通过求解两个系数矩阵都是 M 的方程组来实现, 即每一步需要求解 $Mz = r_{k+1}$ 和 $Mw = Ap_{k+1}$ 来得到 z_{k+1} 和 w_{k+1} . 尽管如此, 若预处理子 $M = LL^T$ 是通过某种矩阵分解 (比如不完全 Cholesky 分解) 得到的, 此时 L 是下三角矩阵, 则 $Mz = r_{k+1}$ 可转化为求解两个三角形方程组 $Ly = r_{k+1}$ 和 $L^T z = y$, 后者的运算量要比前者少得多.

PMINRES 方法的 MATLAB 程序如下:

%PMINRES 方法程序-pminres.m

function [x,iter,time,res,resvec]=pminres1(A,b,x,M1,M2,max_it,tol)

%PMINRES 方法求解对称不定方程组 $Ax=b$, 预处理子 $M=M1*M2$

```

tic; n=length(b);
r=b-A*x; z=M2\ (M1\r); p=z; mr=norm(r);
u=A*p; v=u; w=M2\ (M1\u); iter=1;
while (iter<max_it)
    alpha=(z'*u)/(w'*u); x=x+alpha*p;
    r=r-alpha*u; z1=M2\ (M1\r); v1=A*z1;
    beta=(v1'*z1)/(v'*z); p=z1+beta*p;
    u=v1+beta*u; w=M2\ (M1\u);
    res=norm(r)/mr; resvec(iter)=res;
    if (res<tol), break; end
    z=z1; v=v1; iter=iter+1;
end
time=toc;

```

例 4.9 仍考虑例 4.8 中的线性方程组 $Ax = b$, 取预处理矩阵 $M \approx LL^T$ 为 $A - \mu I$ 的不完全 Cholesky 分解, 其中 $\mu = -1.0$ 是矩阵 A 的最小特征值的一个下界估计. 应用算法 4.11 到该线性方程组上, 第 79 步得到的 $\hat{x} = x_{79}$ 满足

$$\|\hat{x} - x^*\|_2 = 3.2414 \times 10^{-8}.$$

比较极小残量法和预处理极小残量法的数值表现, 迭代过程的收敛轨迹如图 4.10 所示, 其中横坐标为迭代步数 k , 纵坐标为相对残差 $\|r_k\|_2/\|r_0\|_2$, 这里 r_k 是第 k 步得到的残差向量.

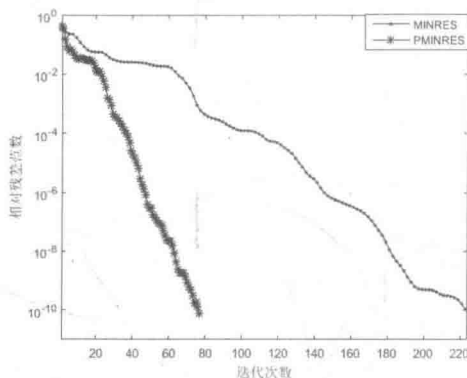


图 4.10 PMINRES 方法与 MINRES 方法的比较

相比于 MINRES 方法, PMINRES 方法尽管迭代次数得到了显著的下降, 从计算时间上看似乎没有优势, 原因可能是后者每一步需要求解两个以预处理子 M 为系数矩阵的线性方程组.

4.3.3 收敛性分析

与共轭梯度法一样, 虽然从理论上讲极小残量法是一种直接方法, 但实际使用时是把它作为迭代法来使用的. 因此, 下面简要说明其收敛特性.

由 Krylov 子空间的性质, $\mathbf{x} \in \mathbf{x}_0 + \mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$ 的充分必要条件是存在 $\varphi \in \mathcal{P}_{k-1}$ 使得 $\mathbf{x} = \mathbf{x}_0 + \varphi(\mathbf{A})\mathbf{r}_0$. 注意到 $\mathbf{x}_k \in \mathbf{x}_0 + \mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$, $\mathbf{x}^* = \mathbf{A}^{-1}\mathbf{b}$, 因此, 有

$$\mathbf{x}^* - \mathbf{x}_k = \mathbf{A}^{-1}\mathbf{b} - \mathbf{x}_0 - \varphi(\mathbf{A})\mathbf{r}_0 = \mathbf{A}^{-1}(\mathbf{I} - \mathbf{A}\varphi(\mathbf{A}))\mathbf{r}_0 = \mathbf{A}^{-1}\psi(\mathbf{A})\mathbf{r}_0,$$

式中: $\psi(t) = 1 - t\varphi(t)$ 满足 $\psi(0) = 1$. 设 \mathbf{A} 的谱分解为 $\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$, 其中 $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n]$ 是正交矩阵, $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$, $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$. 现将 \mathbf{r}_0 按 \mathbf{A} 的特征向量展开, 有

$$\mathbf{r}_0 = \eta_1\mathbf{v}_1 + \eta_2\mathbf{v}_2 + \dots + \eta_n\mathbf{v}_n = \sum_{i=1}^n \eta_i\mathbf{v}_i.$$

于是, 由式 (4.120), 有

$$\begin{aligned} \|\mathbf{r}_k\|_2^2 &= (\mathbf{A}(\mathbf{x}^* - \mathbf{x}_k), \mathbf{A}(\mathbf{x}^* - \mathbf{x}_k)) = (\psi(\mathbf{A})\mathbf{r}_0, \psi(\mathbf{A})\mathbf{r}_0) \\ &= \left(\sum_{i=1}^n \eta_i \psi(\mathbf{A})\mathbf{v}_i, \sum_{j=1}^n \eta_j \psi(\mathbf{A})\mathbf{v}_j \right) \\ &= \left(\sum_{i=1}^n \eta_i \psi(\lambda_i)\mathbf{v}_i, \sum_{j=1}^n \eta_j \psi(\lambda_j)\mathbf{v}_j \right) \\ &= \sum_{i=1}^n \eta_i^2 [\psi(\lambda_i)]^2 \leq \max_{1 \leq i \leq n} [\psi(\lambda_i)]^2 \sum_{i=1}^n \eta_i^2 \\ &= \max_{1 \leq i \leq n} [\psi(\lambda_i)]^2 \|\mathbf{r}_0\|_2^2. \end{aligned}$$

由此可得

$$\|\mathbf{r}_k\|_2 \leq \min_{\psi \in \mathcal{P}_k^0} \max_{\lambda \in \lambda(\mathbf{A})} |\psi(\lambda)| \cdot \|\mathbf{r}_0\|_2, \quad (4.130)$$

式中: 集合 \mathcal{P}_k^0 的定义为

$$\mathcal{P}_k^0 = \{\psi \in \mathcal{P}_k : \psi(0) = 1\}.$$

令

$$a = \min_{\lambda \in \lambda(\mathbf{A})} |\lambda|, \quad b = \max_{\lambda \in \lambda(\mathbf{A})} |\lambda|, \quad (4.131)$$

即 a 和 b 分别是 \mathbf{A} 的特征值的最小模和最大模. 由 \mathbf{A} 是非奇异的假定可知 $0 < a \leq b$. 下面假定 $a < b$, 当 $k = 2m$ 时, 特别取一个多项式

$$\hat{\psi}(t) = C_m \left(\frac{b^2 + a^2 - 2t^2}{b^2 - a^2} \right) / C_m \left(\frac{b^2 + a^2}{b^2 - a^2} \right), \quad (4.132)$$

式中: C_m 为 m 次 Chebyshev 多项式, 即

$$C_m(t) = \frac{(t + \sqrt{t^2 - 1})^m + (t - \sqrt{t^2 - 1})^m}{2},$$

则 $\hat{\psi}(t)$ 是一个 $2m$ 次多项式, 且 $\hat{\psi}(0) = 1$, 从而 $\hat{\psi} \in \mathcal{P}_k^0$. 由 Chebyshev 逼近定理, $\hat{\psi}(t)$ 在区间 $[-b, -a] \cup [a, b]$ 上均有

$$|\hat{\psi}(t)| \leq 1/C_m \left(1 + \frac{2}{\kappa^2 - 1}\right), \quad \forall t \in [-b, -a] \cup [a, b].$$

这样, 由式 (4.130) 即得

$$\frac{\|r_{2m}\|_2}{\|r_0\|_2} \leq \max_{\lambda \in \lambda(A)} |\hat{\psi}(t)| = 1/C_m \left(\frac{b^2 + a^2}{b^2 - a^2}\right) = 1/C_m \left(1 + \frac{2}{\kappa^2 - 1}\right),$$

式中: $\kappa = \|A\|_2 \|A^{-1}\|_2 = b/a$. 记 $\theta(\kappa) = 1 + \frac{2}{\kappa^2 - 1}$, 因 $\|r_{2m+1}\|_2 \leq \|r_{2m}\|_2$, 故有

$$\frac{\|r_k\|_2}{\|r_0\|_2} \leq \frac{1}{C_{\lfloor k/2 \rfloor}(\theta(\kappa))}. \quad (4.133)$$

当然式 (4.133) 的估计是十分粗糙的, 但也说明了极小残量法收敛得快慢与 A 的条件数有关. 事实上, 更精细的分析可以证明, 与共轭梯度法类似, 极小残量法的收敛快慢也与 A 的谱分布有关. 如果 A 的特征值除去少数几个外大多数集中在某数的附近, 则极小残量法将收敛得很快. 因此实际使用时也是将其与预处理技术相结合来使用的.

4.4 SYMMLQ 方法

SYMMLQ 方法与 MINRES 方法类似, 也是求解对称不定方程组的有效算法之一. 本节介绍这一算法的基本思想和导出过程, 并介绍与之关系密切的极小误差法, 以及这两个方法的收敛性及误差估计定理.

4.4.1 SYMMLQ 方法

对于对称不定方程组

$$Ax = b, \quad (4.134)$$

令 $z = x - x_0$, 则式 (4.134) 可改写为

$$Az = r_0, \quad r_0 = b - Ax_0.$$

若

$$S_k = \text{span}\{v_1, v_2, \dots, v_k\}$$

张成一个 k 维子空间. 考虑二次泛函

$$J_k(z_k) = (z - z_k)^T A(z - z_k), \quad (4.135)$$

式中:

$$z_k = V_k y_k, \quad V_k = [v_1, v_2, \dots, v_k].$$

那么二次泛函

$$J_k(z_k) = \tilde{J}_k(y_k) = (A^{-1}r_0 - V_k y_k)^T A (A^{-1}r_0 - V_k y_k)$$

的驻点 y_k 满足

$$V_k^T (r_0 - AV_k y_k) = 0,$$

或者

$$V_k^T AV_k y_k = V_k^T r_0, \quad x_k = x_0 + z_k = x_0 + V_k y_k. \quad (4.136)$$

若取 $v_1 = r_0 / \|r_0\|_2$, 则用对称 Lanczos 算法可得到 k 次 Krylov 子空间

$$\mathcal{K}_k(A, r_0) = \text{span}\{r_0, Ar_0, \dots, A^{k-1}r_0\}$$

的正交规范基 $\{v_i\}_{i=1}^k$.

假设已经得到了一个长度为 k 的 Lanczos 分解 (2.36):

$$AV_k = V_k T_k + \beta_{k+1} v_{k+1} e_k^T, \quad V_k^T V_k = I_k, \quad V_k^T v_{k+1} = 0, \quad (4.137)$$

式中: 三对角矩阵 T_k 通常记为

$$T_k \equiv V_k^T AV_k = \begin{bmatrix} \alpha_1 & \beta_2 & & & \\ \beta_2 & \alpha_2 & \beta_3 & & \\ & \ddots & \ddots & \ddots & \\ & & \beta_{k-1} & \alpha_{k-1} & \beta_k \\ & & & \beta_k & \alpha_k \end{bmatrix}.$$

将上式代入式 (4.136) 则得到 Galerkin 方程组

$$T_k y_k = \beta_1 e_1, \quad x_k = x_0 + V_k y_k, \quad (4.138)$$

式中: $\beta_1 = \|r_0\|_2$.

因为 A 对称不定, 所以三对角矩阵 T_k 也对称不定. 对这样的 T_k 进行 Cholesky 分解可能失效. 因此要寻找 T_k 更为稳定的分解, 即正交三角分解 (LQ 分解):

$$T_k = \tilde{L}_k Q_k, \quad Q_k^T Q_k = I_k, \quad (4.139)$$

式中: \tilde{L}_k 为下三角矩阵.

利用 LQ 分解式 (4.139) 来求解式 (4.138). 和前面的做法一样, 不是直接求出 y_k , 而是寻找求 x_k 递推公式.

因 T_k 是对称三对角矩阵, Q_k 可写为乘积

$$Q_k = G_{k-1,k} \cdots G_{2,3} G_{1,2}, \quad (4.140)$$

式中: $G_{i,i+1}$ ($i = 1, 2, \dots, k-1$) 为 $(i, i+1)$ Givens 矩阵, 它与单位矩阵不同的是其 i 和 $i+1$ 行和列组成的 2 阶矩阵是

$$G_{i,i+1} = \begin{bmatrix} c_i & s_i \\ s_i & -c_i \end{bmatrix}, \quad c_i^2 + s_i^2 = 1.$$

将式 (4.139) 代入式 (4.138), 得

$$\tilde{L}_k(Q_k y_k) = \beta_1 e_1, \quad x_k = x_0 + V_k Q_k^T(Q_k y_k).$$

令

$$\widetilde{W}_k = V_k Q_k^T, \quad \tilde{z}_k = Q_k y_k, \quad (4.141)$$

则得到方程组和迭代向量公式

$$\tilde{L}_k \tilde{z}_k = \beta_1 e_1, \quad x_k = x_0 + \widetilde{W}_k \tilde{z}_k. \quad (4.142)$$

由 Q_k 的特殊形式 (4.140), \widetilde{W}_k 和 \tilde{z}_k 可表示为

$$\widetilde{W}_k = [w_1, w_2, \dots, w_{k-1}, \tilde{w}_k], \quad (4.143)$$

$$\tilde{z}_k = [\zeta_1, \zeta_2, \dots, \zeta_{k-1}, \tilde{\zeta}_k]^T. \quad (4.144)$$

而 LQ 分解式 (4.139) 中的 \tilde{L}_k 和 $G_{k,k+1}$ 的元素 c_k 和 s_k 可计算如下.

由式 (4.139), 公式

$$\tilde{L}_k = T_k Q_k^T = T_k G_{1,2} G_{2,3} \cdots G_{k,k+1}$$

可表示为

$$\tilde{L}_k = \begin{bmatrix} \nu_1 & & & & \\ \delta_2 & \nu_2 & & & \\ \eta_3 & \delta_3 & \nu_3 & & \\ & \ddots & \ddots & \ddots & \\ & & \eta_k & \delta_k & \tilde{\nu}_k \end{bmatrix}. \quad (4.145)$$

现观察 \tilde{L}_{k+1} , 有

$$\begin{aligned} \tilde{L}_{k+1} &= T_{k+1} Q_{k+1}^T = T_{k+1} \left[\begin{array}{c|c} Q_k & \\ \hline & 1 \end{array} \right] G_{k,k+1} \\ &= \left[\begin{array}{cccc|cc} \nu_1 & & & & & \\ \delta_2 & \nu_2 & & & & \\ \eta_3 & \delta_3 & \nu_3 & & & \\ & \ddots & \ddots & \ddots & & \\ & & \eta_k & \delta_k & \tilde{\nu}_k & \beta_{k+1} \\ \hline & & & \eta_{k+1} & \tilde{\delta}_{k+1} & \alpha_{k+1} \end{array} \right] \left[\begin{array}{c|c} I_{k-1} & \\ \hline & c_k \quad s_k \\ & s_k \quad -c_k \end{array} \right], \end{aligned}$$

式中: c_k 和 s_k 由方程组

$$\tilde{\nu}_k c_k + \beta_{k+1} s_k = \nu_k, \quad \tilde{\nu}_k s_k - \beta_{k+1} c_k = 0$$

确定. 由此得 (注意到 $\tilde{\nu}_1 = \alpha_1, \tilde{\delta}_2 = \beta_2$)

$$\nu_k = \sqrt{\tilde{\nu}_k^2 + \beta_{k+1}^2}, \quad c_k = \frac{\tilde{\nu}_k}{\nu_k}, \quad s_k = \frac{\beta_{k+1}}{\nu_k}, \quad (4.146)$$

以及

$$\delta_{k+1} = c_k \tilde{\delta}_{k+1} + s_k \alpha_{k+1}, \quad \tilde{\nu}_{k+1} = s_k \tilde{\delta}_{k+1} - c_k \alpha_{k+1}, \quad (4.147)$$

$$\tilde{\delta}_{k+2} = -c_k \beta_{k+2}, \quad \eta_{k+2} = s_k \beta_{k+2}. \quad (4.148)$$

由式 (4.143) 和式 (4.144), \mathbf{x}_k 在式 (4.142) 中并不能递推地计算. 为此, 引进 \mathbf{L}_k , 它是 $\tilde{\mathbf{L}}_k$ 以 ν_k 替换 $\tilde{\nu}_k$ 而得. 相应地, 定义

$$\mathbf{W}_k = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k], \quad \mathbf{z}_k = [\zeta_1, \zeta_2, \dots, \zeta_k]^T,$$

式中: \mathbf{z}_k 满足三角方程组

$$\mathbf{L}_k \mathbf{z}_k = \beta_1 \mathbf{e}_1. \quad (4.149)$$

由此得 (见 (4.146))

$$\zeta_k = \frac{\tilde{\nu}_k}{\nu_k} \tilde{\zeta}_k = c_k \tilde{\zeta}_k. \quad (4.150)$$

最后, \mathbf{W}_k 的列向量可按下列公式递推计算, 即

$$[\tilde{\mathbf{w}}_k, \mathbf{v}_{k+1}] \begin{bmatrix} c_k & s_k \\ s_k & -c_k \end{bmatrix} = [\mathbf{w}_k, \tilde{\mathbf{w}}_{k+1}], \quad \tilde{\mathbf{w}}_1 = \mathbf{v}_1. \quad (4.151)$$

若 $\tilde{\nu}_k = 0$, 则 $\tilde{\mathbf{L}}_k$ 是奇异的, 故式 (4.142) 不可解. 但由式 (4.146) 可知, 只要 $\beta_{k+1} \neq 0$, \mathbf{L}_k 非奇异, 因而式 (4.149) 可解. 利用 \mathbf{W}_k 和 \mathbf{z}_k 递推地计算 $\tilde{\mathbf{x}}_k$:

$$\tilde{\mathbf{x}}_k = \mathbf{x}_0 + \mathbf{W}_k \mathbf{z}_k = \tilde{\mathbf{x}}_{k-1} + \zeta_k \mathbf{w}_k. \quad (4.152)$$

由式 (4.142), 式 (4.143), 式 (4.144) 和式 (4.152), 有

$$\mathbf{x}_{k+1} = \mathbf{x}_0 + \tilde{\mathbf{W}}_{k+1} \tilde{\mathbf{z}}_{k+1} = \tilde{\mathbf{x}}_k + \tilde{\zeta}_{k+1} \tilde{\mathbf{w}}_{k+1}. \quad (4.153)$$

式 (4.150) 只是给出了 ζ_k 与 $\tilde{\zeta}_k$ 之间的关系. 还需要 ζ_k 或 $\tilde{\zeta}_k$ ($k=1, 2, \dots$) 的计算公式. 由式 (4.149), 得

$$\begin{cases} \nu_1 \zeta_1 = \beta_1 \implies \zeta_1 = \frac{\beta_1}{\nu_1}, \\ \delta_2 \zeta_1 + \nu_2 \zeta_2 = 0 \implies \zeta_2 = -\frac{\delta_2 \zeta_1}{\nu_2}, \\ \eta_k \zeta_{k-2} + \delta_k \zeta_{k-1} + \nu_k \zeta_k = 0 \implies \zeta_k = -\frac{\eta_k \zeta_{k-2} + \delta_k \zeta_{k-1}}{\nu_k}, \quad k \geq 3. \end{cases} \quad (4.154)$$

求解式 (4.149) 有比式 (4.142) 更好的数值性态. 由式 (4.152) 可求得 $\{\tilde{\mathbf{x}}_k\}$, 若需要 $\{\mathbf{x}_k\}$ 可按式 (4.153) 得到. 所描述的算法称为 SYMMLQ 方法.

算法 4.12 (SYMMMLQ 方法) 给定 n 阶非奇异的实对称矩阵 A , n 维向量 b 和允许误差 $\varepsilon > 0$. 本算法计算向量 x_k , 使得 $\|r_k\|_2/\|r_0\|_2 \leq \varepsilon$, 其中 $r_k = b - Ax_k$.

选取 x_0 ; $r_0 = b - Ax_0$; $\beta_1 = \|r_0\|_2$; $v_1 = r_0/\beta_1$;

$\alpha_1 = v_1^T A v_1$; $\tilde{v}_2 = A v_1 - \alpha_1 v_1$; $\beta_2 = \|\tilde{v}_2\|_2$;

$v_2 = \tilde{v}_2/\beta_2$; $\tilde{w}_1 = v_1$;

$\tilde{\nu}_1 = \alpha_1$; $\tilde{\delta}_2 = \beta_2$; $\zeta_0 = 0$; $\eta_2 = 0$;

$\nu_1 = \sqrt{\alpha_1^2 + \beta_2^2}$; $c_1 = \tilde{\nu}_1/\nu_1$; $s_1 = \beta_2/\nu_1$;

由 $[\tilde{w}_1, v_2] \begin{bmatrix} c_1 & s_1 \\ s_1 & -c_1 \end{bmatrix} = [w_1, \tilde{w}_2]$ 确定 w_1 和 \tilde{w}_2 ;

即 $w_1 = c_1 \tilde{w}_1 + s_1 v_2$; $\tilde{w}_2 = s_1 \tilde{w}_1 - c_1 v_2$;

$\zeta_1 = \frac{\beta_1}{\nu_1}$; $\tilde{x}_1 = x_0 + \zeta_1 w_1$; $k = 1$;

while ($\|r_k\|_2/\|r_0\|_2 > \varepsilon$)

$\alpha_{k+1} = v_{k+1}^T A v_{k+1}$;

$\tilde{v}_{k+2} = A v_{k+1} - \alpha_{k+1} v_{k+1} - \beta_{k+1} v_k$;

$\beta_{k+2} = \|\tilde{v}_{k+2}\|_2$; $v_{k+2} = \tilde{v}_{k+2}/\beta_{k+2}$;

$\tilde{\nu}_{k+1} = s_k \tilde{\delta}_{k+1} - c_k \alpha_{k+1}$; $\nu_{k+1} = \sqrt{\tilde{\nu}_{k+1}^2 + \beta_{k+2}^2}$;

$c_{k+1} = \frac{\tilde{\nu}_{k+1}}{\nu_{k+1}}$; $s_{k+1} = \frac{\beta_{k+2}}{\nu_{k+1}}$;

$\tilde{\delta}_{k+1} = c_k \tilde{\delta}_{k+1} + s_k \alpha_{k+1}$;

$\tilde{\delta}_{k+2} = -c_k \beta_{k+2}$; $\eta_{k+2} = s_k \beta_{k+2}$;

$\zeta_{k+1} = -\frac{\eta_{k+1} \zeta_{k-1} + \delta_{k+1} \zeta_k}{\nu_{k+1}}$; $\tilde{\zeta}_{k+1} = \zeta_{k+1}/c_{k+1}$;

由 $[\tilde{w}_{k+1}, v_{k+2}] \begin{bmatrix} c_{k+1} & s_{k+1} \\ s_{k+1} & -c_{k+1} \end{bmatrix} = [w_{k+1}, \tilde{w}_{k+2}]$ 确定 w_{k+1} 和 \tilde{w}_{k+2} ;

即 $w_{k+1} = c_{k+1} \tilde{w}_{k+1} + s_{k+1} v_{k+2}$; $\tilde{w}_{k+2} = s_{k+1} \tilde{w}_{k+1} - c_{k+1} v_{k+2}$;

$x_{k+1} = \tilde{x}_k + \tilde{\zeta}_{k+1} \tilde{w}_{k+1}$; $\tilde{x}_{k+1} = \tilde{x}_k + \zeta_{k+1} w_{k+1}$;

$r_{k+1} = b - A x_{k+1}$;

$k = k + 1$;

end

例 4.10 仍考虑例 4.8 中的线性方程组 $Ax = b$, 其中矩阵 A 来自

<http://www.cise.ufl.edu/research/sparse/matrices/PARSEC/SiNa.html>,

它是一个 5743 阶的实对称不定稀疏矩阵, 右端项 $b = Ae$, 其中 $e = (1, 1, \dots, 1)^T$. 显然, 该方程组的真解为 $x^* = (1, 1, \dots, 1)^T$. 将 SYMMMLQ 方法应用到该线性方程组上, 迭代在 224 步后收敛 ($\varepsilon = 10^{-10}$), 计算得到的近似解 \hat{x} 和真解 x^* 之间的绝对值误差为

$$\|\hat{x} - x^*\|_2 = 2.8998 \times 10^{-8},$$

计算解 \hat{x} 的残量满足

$$\|b - A\hat{x}\|_2 = 2.1331 \times 10^{-8}.$$

将它与 MINRES 方法和 GMRES 方法进行了对比, 发现 SYMMLQ 方法与 MINRES 方法对本例有几乎完全等效的数值表现, 而 GMRES 方法则要略好一些. 迭代过程的收敛轨迹如图 4.11 所示, 其中横坐标为迭代步数 k , 纵坐标为相对残差 $\|r_k\|_2/\|r_0\|_2$, 这里 r_k 是第 k 步得到的残差向量.

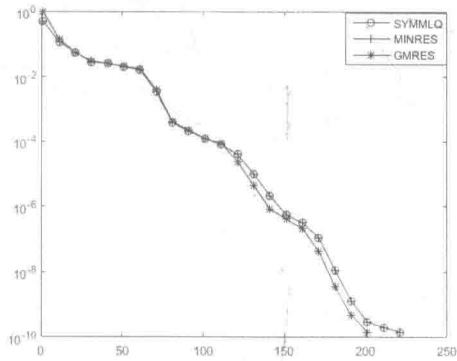


图 4.11 SYMMLQ 方法的收敛特性

4.4.2 收敛性分析

对于方程组

$$Az = r_0, \quad r_0 = b - Ax_0, \quad z = x - x_0,$$

考虑二次泛函

$$J_k^{(k)}(z_k) = (z - z_k)^T A^k (z - z_k), \tag{4.155}$$

式中: k 为非负整数; $z_k = V_k y_k$, $V_k = [v_1, v_2, \dots, v_k]$.

当 $k = 2$ 时, 式 (4.155) 中的泛函

$$\begin{aligned} J_k^{(2)}(z_k) &= (z - z_k)^T A^2 (z - z_k) = \|A(z - z_k)\|_2^2 \\ &= \|r_0 - A(x_k - x_0)\|_2^2 = \|b - Ax_k\|_2^2 = \|r_k\|_2^2 \end{aligned}$$

是正定二次泛函. 极小化这一正定二次泛函就是 4.3 节中极小残量法 (MINRES) 的基本思想.

当 $k = 1$ 时, 式 (4.155) 中的泛函为

$$J_k^{(1)}(z_k) = (z - z_k)^T A (z - z_k) = (x - x_0 - V_k y_k)^T A (x - x_0 - V_k y_k).$$

求上述泛函稳定点 y_k :

$$V_k^T A (x - x_0 - V_k y_k) = V_k^T (r_0 - AV_k y_k) = 0$$

是本节 SYMMLQ 方法的出发点.

若在式 (4.155) 中取 $k = 0$, 则有

$$J_k^{(0)}(z_k) = (z - z_k)^T(z - z_k) = \|x - x_k\|_2^2 = \|\varepsilon_k\|_2^2.$$

这也是一个正定二次泛函. 因而使 $J_k^{(0)}$ 极小化就是使误差 $\varepsilon_k = x - x_k$ 的范数极小化.

为了使极小误差法可行, 取迭代所处的仿射子空间为

$$x_0 + \mathcal{K}_k(A, Ar_0),$$

也就是说

$$z_k = x_k - x_0 = AV_k y_k.$$

由此, y_k 是最小二乘问题

$$\min \|x - x_0 - AV_k y_k\|_2 = \|A^{-1}r_0 - AV_k y_k\|_2$$

的解. 由法方程可知 y_k 满足线性方程组

$$V_k^T A^2 V_k y_k = V_k^T r_0 = \beta_1 e_1, \quad x_k = x_0 + AV_k y_k. \quad (4.156)$$

由对称 Lanczos 分解式 (4.137) 及 LQ 分解式 (4.139), 得

$$\begin{aligned} V_k^T A^2 V_k &= (AV_k)^T (AV_k) \\ &= (V_k T_k + \beta_{k+1} v_{k+1} e_k^T)^T (V_k T_k + \beta_{k+1} v_{k+1} e_k^T) \\ &= T_k^2 + \beta_{k+1}^2 e_k e_k^T = (\tilde{L}_k Q_k)(\tilde{L}_k Q_k)^T + \beta_{k+1}^2 e_k e_k^T \\ &= \tilde{L}_k \tilde{L}_k^T + \beta_{k+1}^2 e_k e_k^T = L_k L_k^T, \end{aligned}$$

最后一个等式成立是因为 $\tilde{L}_k \tilde{L}_k^T$ 与 $L_k L_k^T$ 只是 (k, k) 处的那个元素不同, 前者是 $\eta_k^2 + \delta_k^2 + \tilde{\nu}_k^2$, 而后者为 $\eta_k^2 + \delta_k^2 + \nu_k^2$. 注意到 $\nu_k^2 = \tilde{\nu}_k^2 + \beta_{k+1}^2$, 即可得等式成立. 这样, 直接从 T_k 的 LQ 分解得到了 $V_k^T A^2 V_k$ 的 Cholesky 分解. 代入式 (4.156), 需要求解如下线性方程组以得到 y_k 和 x_k :

$$L_k L_k^T y_k = \beta_1 e_1, \quad x_k = x_0 + AV_k y_k. \quad (4.157)$$

由上述方式得到近似解序列 $\{x_k\}$ 的方法通常称为极小误差法. 下面证明这个极小误差解序列 $\{x_k\}$ 就是 SYMMLQ 方法中的 $\{\tilde{x}_k\}$ (见式 (4.152)).

定理 4.9 SYMMLQ 方法中的序列 $\{\tilde{x}_k\}$ 是在 Krylov 子空间 $\mathcal{K}_k(A, Ar_0)$ 的极小误差逼近.

证明 由式 (4.157) 可知极小误差解 x_k 可表示为

$$x_k = x_0 + \beta_1 AV_k L_k^{-T} L_k^{-1} e_1. \quad (4.158)$$

由式 (4.149) 和式 (4.152) 可知, $\tilde{\mathbf{x}}_k$ 可表示为

$$\tilde{\mathbf{x}}_k = \mathbf{x}_0 + \beta_1 \mathbf{W}_k \mathbf{L}_k^{-1} \mathbf{e}_1. \quad (4.159)$$

比较式 (4.158) 和式 (4.159) 右端, 只需证明

$$\mathbf{W}_k = \mathbf{A} \mathbf{V}_k \mathbf{L}_k^{-\mathrm{T}}, \quad (4.160)$$

即完成了定理的证明. 由式 (4.139) 和式 (4.141) 可知

$$\tilde{\mathbf{W}}_{k+1} \tilde{\mathbf{L}}_{k+1}^{\mathrm{T}} = \mathbf{V}_{k+1} (\tilde{\mathbf{L}}_{k+1} \mathbf{Q}_{k+1})^{\mathrm{T}} = \mathbf{V}_{k+1} \mathbf{T}_{k+1}. \quad (4.161)$$

再由式 (4.137), 得

$$\tilde{\mathbf{W}}_{k+1} = (\mathbf{A} \mathbf{V}_{k+1} - \beta_{k+2} \mathbf{v}_{k+2} \mathbf{e}_{k+1}^{\mathrm{T}}) \tilde{\mathbf{L}}_{k+1}^{-\mathrm{T}}. \quad (4.162)$$

比较式 (4.162) 前 k 列, 并注意到 $\tilde{\mathbf{L}}_{k+1}^{-\mathrm{T}}$ 是上三角矩阵, 可得式 (4.160). 证毕. \square

下面分析极小误差法的误差估计. 令

$$\tilde{\mathcal{P}}_{k-1}^0(t) = 1 - t^2 \mathcal{P}_{k-1}(t)$$

表示次数不超过 $k+1$, 在 $t=0$ 时值等于 1 且一阶导数等于 0 的多项式集合. 因 $\mathbf{x}_k - \mathbf{x}_0 \in \mathcal{K}_k(\mathbf{A}, \mathbf{A} \mathbf{r}_0)$, 故存在 $p_{k-1} \in \mathcal{P}_{k-1}$, 使得 $\mathbf{x}_k - \mathbf{x}_0 = p_{k-1}(\mathbf{A}) \mathbf{A} \mathbf{r}_0$. 所以误差 $\varepsilon_k = \mathbf{x}^* - \mathbf{x}_k$ 满足

$$\begin{aligned} \|\varepsilon_k\|_2^2 &= \min_{p_{k-1} \in \mathcal{P}_{k-1}} (\varepsilon_0 - p_{k-1}(\mathbf{A}) \mathbf{A} \mathbf{r}_0, \varepsilon_0 - p_{k-1}(\mathbf{A}) \mathbf{A} \mathbf{r}_0) \\ &= \min_{p_{k-1} \in \mathcal{P}_{k-1}} (\varepsilon_0 - p_{k-1}(\mathbf{A}) \mathbf{A}^2 \varepsilon_0, \varepsilon_0 - p_{k-1}(\mathbf{A}) \mathbf{A}^2 \varepsilon_0) \\ &= \min_{p_{k-1} \in \mathcal{P}_{k-1}} ([\mathbf{I} - \mathbf{A}^2 p_{k-1}(\mathbf{A})] \varepsilon_0, [\mathbf{I} - \mathbf{A}^2 p_{k-1}(\mathbf{A})] \varepsilon_0) \\ &= \min_{\tilde{p}_{k-1} \in \tilde{\mathcal{P}}_{k-1}^0} (\tilde{p}_{k-1}(\mathbf{A}) \varepsilon_0, \tilde{p}_{k-1}(\mathbf{A}) \varepsilon_0) \\ &\leq \min_{\tilde{p}_{k-1} \in \tilde{\mathcal{P}}_{k-1}^0} \max_{\lambda \in \sigma(\mathbf{A})} \tilde{p}_{k-1}^2(\lambda) \|\varepsilon_0\|_2^2. \end{aligned} \quad (4.163)$$

因为 \mathbf{A} 是对称不定的, 所以 $\sigma(\mathbf{A})$ 包含正的和负的特征值. 有下面的定理.

定理 4.10 若 \mathbf{A} 的谱 $\sigma(\mathbf{A}) \subset [-b, -a] \cup [a, b]$, 其中 $0 < a < b$. 则对极小误差迭代法成立

$$\|\varepsilon_k\|_2 \leq \frac{2\tilde{r}^{\tilde{k}}}{1 + \tilde{r}^{2\tilde{k}}} \|\varepsilon_0\|_2, \quad (4.164)$$

式中:

$$\tilde{r} = \frac{b-a}{b+a} = \frac{b/a-1}{b/a+1}, \quad \tilde{k} = \left\lceil \frac{k+1}{2} \right\rceil.$$

证明 取多项式

$$p_{k-1}(t) = \frac{C_{\tilde{k}}(v(t))}{C_{\tilde{k}}(v(0))},$$

式中: $C_{\tilde{k}}(z)$ 为 $\tilde{k} = [(k+1)/2]$ 次 Chebyshev 多项式, 且

$$v(t) = 1 - \frac{2(t^2 - a^2)}{b^2 - a^2},$$

它将区间 $[-b, -a] \cup [a, b]$ 映射到 $[-1, 1]$.

显然, $p_{k-1}(t)$ 的次数不超过 $k+1$, 它是偶函数, 且 $p_{k-1}(0) = 1$. 故 $p_{k-1}(t) \in \tilde{P}_{k-1}^0$. 由式 (4.163), 有

$$\|\varepsilon_k\|_2 \leq \max_{t \in [-b, -a] \cup [a, b]} \left| \frac{C_{\tilde{k}}(v(t))}{C_{\tilde{k}}(v(0))} \right| \|\varepsilon_0\|_2 = \frac{1}{C_{\tilde{k}}(v(0))} \|\varepsilon_0\|_2.$$

不难推得

$$\begin{aligned} C_{\tilde{k}}(v(0)) &= C_{\tilde{k}}\left(\frac{b^2 + a^2}{b^2 - a^2}\right) \\ &= \frac{1}{2} \left\{ \left[\frac{b^2 + a^2}{b^2 - a^2} + \sqrt{\left(\frac{b^2 + a^2}{b^2 - a^2}\right)^2 - 1} \right]^{\tilde{k}} + \left[\frac{b^2 + a^2}{b^2 - a^2} - \sqrt{\left(\frac{b^2 + a^2}{b^2 - a^2}\right)^2 - 1} \right]^{-\tilde{k}} \right\} \\ &= \frac{1}{2} \left[\left(\frac{b+a}{b-a}\right)^{\tilde{k}} + \left(\frac{b-a}{b+a}\right)^{-\tilde{k}} \right] = \frac{1}{2} \left[\left(\frac{1}{\tilde{r}}\right)^{\tilde{k}} + \tilde{r}^{\tilde{k}} \right] = \frac{1 + \tilde{r}^{2\tilde{k}}}{2\tilde{r}^{\tilde{k}}}. \end{aligned}$$

证毕. □

4.5 拟极小残量法

在 4.2 节中, 为了解决广义极小残量法的存储问题而采用了重新开始的方法, 进而导出了实际应用中常用的重开始广义极小残量法 GMRES(m). 本节介绍另一种解决存储问题的方法—拟极小残量法, 在目前的文献中常简称为 QMRES 方法 (Quasi-Minimal RESidual Approach).

GMRES 方法的存储量主要来自计算 Krylov 子空间 $\mathcal{K}_k(\mathbf{A}, \mathbf{b})$ 的正交基 $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$. 由于要求这些 \mathbf{v}_i ($i = 1, 2, \dots, k$) 相互正交, 故必须在计算中将这组向量都存储起来. 因此要使其存储量降低, 一个自然的想法是放弃正交性要求, 而选择其他适当的基向量.

类比于 MINRES 方法, 自然希望计算 Krylov 子空间 $\mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$ 的一组基向量 $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$, 使得

$$\mathcal{K}_k(\mathbf{A}, \mathbf{r}_0) = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_k\}, \quad \mathbf{v}_1 = \mathbf{r}_0 / \|\mathbf{r}_0\|_2, \quad (4.165)$$

$$\mathbf{A}\mathbf{V}_k = \mathbf{V}_{k+1}\tilde{\mathbf{T}}_k, \quad (4.166)$$

式中: $\mathbf{V}_{k+1} = [\mathbf{V}_k, \mathbf{v}_{k+1}] = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k, \mathbf{v}_{k+1}]$ 列满秩, 而 $\tilde{\mathbf{T}}_k$ 是 $(k+1) \times k$ 阶三对角矩阵.

如果能够找到一种方法来实现式 (4.165) 和式 (4.166) 中的分解, 则由于这样的基满足等式 $\mathbf{A}\mathbf{V}_k = \mathbf{V}_{k+1}\tilde{\mathbf{T}}_k$, 向量 \mathbf{v}_i 就可由三项递推公式来确定, 从而所需的存储量只有 $O(n)$ 而并非 GMRES 方法的 $O(kn)$. 然而, 这样一来, 所要解决的极小化问题 (4.52) 就转化为求 $\mathbf{z}_k \in \mathbb{R}^k$, 使得

$$\|\mathbf{r}_k\|_2 = \|\mathbf{V}_{k+1}(\beta\mathbf{e}_1 - \tilde{\mathbf{T}}_k\mathbf{z}_k)\|_2 = \min, \quad (4.167)$$

式中: $\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{V}_k\mathbf{z}_k$, $\beta = \|\mathbf{r}_0\|_2$. 但由于此时的 \mathbf{V}_{k+1} 列向量并不正交, 故求解这一极小化问题就有很大的困难. 注意到

$$\|\mathbf{r}_k\|_2 = \|\mathbf{V}_{k+1}(\beta\mathbf{e}_1 - \tilde{\mathbf{T}}_k\mathbf{z}_k)\|_2 \leq \|\mathbf{V}_{k+1}\|_2 \|\beta\mathbf{e}_1 - \tilde{\mathbf{T}}_k\mathbf{z}_k\|_2,$$

如果 $\|\mathbf{V}_{k+1}\|_2$ 不是特别大, 则只需求得最小二乘问题

$$\min \{ \|\beta\mathbf{e}_1 - \tilde{\mathbf{T}}_k\mathbf{z}\|_2 : \mathbf{z} \in \mathbb{R}^k \} \quad (4.168)$$

的解 \mathbf{z}_k . 虽然 $\|\mathbf{r}_k\|_2$ 并没有达到最小, 但也会相对很小, 这就是拟极小化方法的基本想法. 总结起来, 主要有两步:

第 1 步, 计算 $\mathcal{K}_k(\mathbf{A}, \mathbf{b})$ 之满足条件 (4.165) 和 (4.166) 的一组基.

第 2 步, 求解最小二乘问题 (4.168).

关键是第 1 步如何实现, 第 2 步可用与 4.3 节中完全一样的方法来求解. 这里用非对称 Lanczos 方法实现 $\mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$ 基向量的计算问题.

4.5.1 非对称 Lanczos 方法

非对称 Lanczos 方法是对称 Lanczos 方法的一种自然的推广, 它是同时计算 $\mathcal{K}_k(\mathbf{A}, \mathbf{v}_1)$ 和 $\mathcal{K}_k(\mathbf{A}^T, \mathbf{w}_1)$ 的基向量 $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$ 和 $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k$, 使得

$$\mathbf{w}_i^T \mathbf{v}_j = \begin{cases} 0, & i \neq j, \\ 1, & i = j. \end{cases} \quad (4.169)$$

式 (4.169) 称为双正交性条件.

设矩阵 \mathbf{A} 和向量 \mathbf{v}_1 已经给定, 并且假定 $\|\mathbf{v}_1\|_2 = 1$. 首先需选择一个 $\mathbf{w}_1 \in \mathbb{R}^n$, 使得 $\mathbf{v}_1^T \mathbf{w}_1 = 1$. 如果没有什么特别的信息可以利用的话, 可以简单地选择 $\mathbf{w}_1 = \mathbf{v}_1$. 然后类比于对称 Lanczos 方法的计算过程, 也采用如下的三项递推方式, 即

$$\begin{cases} \beta_i \mathbf{v}_{i+1} = \mathbf{A}\mathbf{v}_i - \alpha_i \mathbf{v}_i - \gamma_{i-1} \mathbf{v}_{i-1}, \\ \gamma_i \mathbf{w}_{i+1} = \mathbf{A}^T \mathbf{w}_i - \alpha_i \mathbf{w}_i - \beta_{i-1} \mathbf{w}_{i-1}. \end{cases} \quad (4.170)$$

依次确定 $\alpha_i, \beta_i, \gamma_i, \mathbf{v}_{i+1}$ 和 \mathbf{w}_{i+1} . 这里假定 $\gamma_0 \mathbf{v}_0 = \beta_0 \mathbf{w}_0 = \mathbf{0}$, 即当 $i = 1$ 时, 有

$$\beta_1 \mathbf{v}_2 = \mathbf{A}\mathbf{v}_1 - \alpha_1 \mathbf{v}_1,$$

$$\gamma_1 \mathbf{w}_2 = \mathbf{A}^T \mathbf{w}_1 - \alpha_1 \mathbf{w}_1.$$

由于希望计算得到的 \mathbf{v}_2 和 \mathbf{w}_2 满足 $\mathbf{v}_2 \perp \mathbf{w}_1$, $\mathbf{w}_2 \perp \mathbf{v}_1$, 故在上面的第 1 式两边左乘 \mathbf{w}_1^T , 得

$$\alpha_1 = \mathbf{w}_1^T \mathbf{A} \mathbf{v}_1. \quad (4.171)$$

此处利用了 $\mathbf{w}_1^T \mathbf{v}_1 = 1$. 当然, 也可以在前面的第 2 式两边左乘 \mathbf{v}_1^T , 得到 α_1 的表达式 (4.171). 一旦 α_1 确定, 便可计算

$$\begin{aligned} \beta_1 \mathbf{v}_2 &= \tilde{\mathbf{v}}_2 = \mathbf{A} \mathbf{v}_1 - \alpha_1 \mathbf{v}_1, \\ \gamma_1 \mathbf{w}_2 &= \tilde{\mathbf{w}}_2 = \mathbf{A}^T \mathbf{w}_1 - \alpha_1 \mathbf{w}_1. \end{aligned}$$

由于要求 $\mathbf{v}_2^T \mathbf{w}_2 = 1$, 故必有

$$\beta_1 \gamma_1 = \tilde{\mathbf{v}}_2^T \tilde{\mathbf{w}}_2 = \omega_1.$$

如果 $\omega_1 = 0$, 则只好结束; 否则可选择两个数 β_1 和 γ_1 使上式成立. 可选择

$$\beta_1 = \sqrt{|\omega_1|}, \quad \gamma_1 = \omega_1 / \beta_1. \quad (4.172)$$

这样一来, 所需的向量 \mathbf{v}_2 和 \mathbf{w}_2 即为

$$\mathbf{v}_2 = \tilde{\mathbf{v}}_2 / \beta_1, \quad \mathbf{w}_2 = \tilde{\mathbf{w}}_2 / \gamma_1. \quad (4.173)$$

假定已经确定了

$$\begin{aligned} &\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_i; \quad \mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_i; \\ &\alpha_1, \alpha_2, \dots, \alpha_{i-1}; \quad \gamma_1, \gamma_2, \dots, \gamma_{i-1}; \quad \beta_1, \beta_2, \dots, \beta_{i-1}, \end{aligned}$$

下一步利用 $\mathbf{w}_i^T \mathbf{v}_{i+1} = \mathbf{w}_i^T \mathbf{v}_{i-1} = 0$ 和 $\mathbf{w}_i^T \mathbf{v}_i = 1$, 从式 (4.170) 的第 1 式可导出

$$\alpha_i = \mathbf{w}_i^T \mathbf{A} \mathbf{v}_i. \quad (4.174)$$

然后, 计算

$$\begin{cases} \tilde{\mathbf{v}}_{i+1} = \mathbf{A} \mathbf{v}_i - \alpha_i \mathbf{v}_i - \gamma_{i-1} \mathbf{v}_{i-1}, \\ \tilde{\mathbf{w}}_{i+1} = \mathbf{A}^T \mathbf{w}_i - \alpha_i \mathbf{w}_i - \beta_{i-1} \mathbf{w}_{i-1}, \\ \omega_i = \tilde{\mathbf{v}}_{i+1}^T \tilde{\mathbf{w}}_{i+1}. \end{cases} \quad (4.175)$$

若 $\omega_i = 0$, 则结束; 否则, 计算

$$\beta_i = \sqrt{|\omega_i|}, \quad \gamma_i = \omega_i / \beta_i, \quad \mathbf{v}_{i+1} = \tilde{\mathbf{v}}_{i+1} / \beta_i, \quad \mathbf{w}_{i+1} = \tilde{\mathbf{w}}_{i+1} / \gamma_i. \quad (4.176)$$

假定上述过程共进行了 k 步, 现在令

$$\mathbf{V}_k = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k], \quad \mathbf{W}_k = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k],$$

$$T_k = \begin{bmatrix} \alpha_1 & \gamma_1 & & & \\ \beta_1 & \alpha_2 & \gamma_2 & & \\ & \beta_2 & \ddots & \ddots & \\ & & \ddots & \ddots & \gamma_{k-1} \\ & & & \beta_{k-1} & \alpha_k \end{bmatrix}. \quad (4.177)$$

则易从上面的推导过程归纳地证明关于非对称 Lanczos 迭代的如下基本性质:

- (1) $AV_k = V_k T_k + \beta_k v_{k+1} e_k^T$.
- (2) $A^T W_k = W_k T_k^T + \gamma_k w_{k+1} e_k^T$.
- (3) $V_k^T W_k = I_k$, $v_{k+1}^T W_k = 0$, $w_{k+1}^T V_k = 0$, $w_{k+1}^T v_{k+1} = 1$.
- (4) $\mathcal{K}_k(A, v_1) = \text{span}\{v_1, v_2, \dots, v_k\}$.
- (5) $\mathcal{K}_k(A^T, w_1) = \text{span}\{w_1, w_2, \dots, w_k\}$.

其实 (1) 和 (2) 就是前面算法的矩阵表示. (4) 和 (5) 是前三条的直接推论. 例如 (4), 由 (1) 归纳地可证 $v_i \in \mathcal{K}_k(A, v_1)$, $i = 1, 2, \dots, k$, 而 (3) 又蕴涵着 V_k 是列满秩的 (即 v_1, v_2, \dots, v_k 是线性无关的), 从而必有 (4) 成立. 因此关键是 (3) 的证明.

下面用数学归纳法证明性质 (3). 首先由初值 w_1 的选择知, 自然有 $w_1^T v_1 = 1$ 成立. 现在假设已经对自然数 i 证明了向量组 v_1, v_2, \dots, v_i 和 w_1, w_2, \dots, w_i 满足双正交性条件, 再来考虑 $i+1$ 的情形. 注意这里假定 $i+1 \leq k+1$.

首先由归纳法假设和 α_i 的定义, 有

$$\begin{aligned} v_{i+1}^T w_i &= \frac{1}{\beta_i} (A v_i - \alpha_i v_i - \gamma_{i-1} v_{i-1})^T w_i \\ &= \frac{1}{\beta_i} (v_i^T A^T w_i - \alpha_i) = 0. \end{aligned}$$

对任意的 $l < i$, 有

$$\begin{aligned} v_{i+1}^T w_l &= \frac{1}{\beta_i} (v_i^T A^T w_l - \alpha_i v_i^T w_l - \gamma_{i-1} v_{i-1}^T w_l) \\ &= \frac{1}{\beta_i} [v_i^T (\alpha_l w_l + \beta_{l-1} w_{l-1} + \gamma_l w_{l+1}) - \gamma_{i-1} v_{i-1}^T w_l]. \end{aligned}$$

当 $l < i-1$ 时, 由归纳法假设立即由上式知 $v_{i+1}^T w_l = 0$. 而当 $l = i-1$ 时, 直接计算有

$$v_{i+1}^T w_{i-1} = \frac{1}{\beta_i} (\alpha_{i-1} v_i^T w_{i-1} + \beta_{i-2} v_i^T w_{i-2} + \gamma_{i-1} v_i^T w_i - \gamma_{i-1} v_{i-1}^T w_{i-1}) = 0,$$

最后的等式用到了归纳法假设. 至于 $w_{i+1}^T v_{i+1} = 1$ 可由计算过程立即得到. 这表明 (3) 对 $i+1$ 也成立. 因此由归纳法原理知 (3) 成立.

综合上面的讨论, 可得下面的非对称 Lanczos 方法.

算法 4.13 (非对称 Lanczos 方法) 给定矩阵 $A \in \mathbb{R}^{n \times n}$, 向量 $v_1 \in \mathbb{R}^n$ ($\|v_1\|_2 = 1$) 和正整数 k . 本算法计算形如 (4.177) 的三个矩阵 V_k, W_k 和 T_k , 以及向量 v_{k+1}, w_{k+1} 和数 β_k, γ_k , 满足前面所述的 5 条性质.

```

选择向量  $w_1 \in \mathbb{R}^n$  满足  $w_1^T v_1 = 1$ ;
 $v = Av_1$ ;  $w = A^T w_1$ ;
for  $i = 1 : k$ 
     $\alpha_i = w_i^T v$ ;  $v = v - \alpha_i v_i$ ;  $w = w - \alpha_i w_i$ ;
    if  $\|v\|_2 = 0$  或者  $\|w\|_2 = 0$ 
        stop
    else
         $\omega_i = v^T w$ ;
    end
    if  $\omega_i = 0$ 
        stop
    else
         $\beta_i = \sqrt{|\omega_i|}$ ;  $\gamma_i = \omega_i / \beta_i$ ;
         $v_{i+1} = v / \beta_i$ ;  $w_{i+1} = w / \gamma_i$ ;
    end
     $v = Av_{i+1} - \gamma_i v_i$ ;  $w = A^T w_{i+1} - \beta_i w_i$ ;
end
```

注 4.5 非对称 Lanczos 迭代的基本性质 (1) ~ (3) 蕴涵着

$$T_k = W_k^T A V_k.$$

这表明 T_k 正好是 A 沿着 $\mathcal{K}_k(A^T, w_1)^\perp$ 到 $\mathcal{K}_k(A, v_1)$ 上的投影. 当然, T_k^T 也是 A 沿着 $\mathcal{K}_k(A, v_1)^\perp$ 到 $\mathcal{K}_k(A^T, w_1)$ 上的投影. 有关投影的几何解释见注 4.6.

注 4.6 设 $X, Y \in \mathbb{R}^{n \times k}$ 满足 $Y^T X = I_k$, 并记

$$\mathcal{X} = \mathcal{R}(X), \quad \mathcal{Y} = \mathcal{R}(Y), \quad P = XY^T,$$

则容易验证

$$\mathcal{R}(P) = \mathcal{X}, \quad \mathcal{N}(P) = \mathcal{Y}^\perp, \quad P^2 = P,$$

这表明 P 是沿 \mathcal{Y}^\perp 到 \mathcal{X} 上的投影算子.

设 $A \in \mathbb{R}^{n \times n}$, 将 A 视作 \mathbb{R}^n 到 \mathbb{R}^n 的线性算子, 并记其在 \mathcal{X} 上的限制为 $A|_{\mathcal{X}}$, 则 $P \circ A|_{\mathcal{X}}$ 就是从 \mathcal{X} 到 \mathcal{X} 的线性算子. 由于

$$P \circ A|_{\mathcal{X}}(Xv) = XY^T AXv, \quad \forall v \in \mathbb{R}^k,$$

故 $B = Y^T AX$ 正好是算子 $P \circ A|_{\mathcal{X}}$ 在给定基 X 之下的矩阵表示. 因此, 通常就称 B 是 A 沿 \mathcal{Y}^\perp 到 \mathcal{X} 上的投影.

从实际应用的角度来看, 非对称 Lanczos 方法比 Arnoldi 方法有着很大的优势, 这主要体现在它仅需存储六个 n 维向量 (因为 $k \ll n$, 所以 T_k 的存储量是微不足道的), 并不随着 k 的增加而增加.

但是, 从另一个角度来看, 非对称 Lanczos 算法会发生中断, 即会出现 \tilde{v}_{i+1} 和 \tilde{w}_{i+1} 均不为零, 而 $\tilde{w}_{i+1}^T \tilde{v}_{i+1} = 0$ 的情形. 虽然在实际计算时出现 $\tilde{w}_{i+1}^T \tilde{v}_{i+1} = 0$ 的概率很小, 但 $|\tilde{w}_{i+1}^T \tilde{v}_{i+1}|$ 很小的情况还是经常会遇到的, 此时 $|\gamma_i|$ 和 $|\beta_i|$ 之中有一个就会变得很小, 从而导致在计算过程中引进较大的误差. 解决这一问题的一种方法就是采用 Look-ahead 技术. Look-ahead 方式是基于对算法 4.13 的仔细观察而得到的. 在遇到中断的时候, 虽然不能给出 v_{i+1} 和 w_{i+1} , 但却常常可以给出 v_{i+2} 和 w_{i+2} . 这就使得 Lanczos 方法可以继续下去. Look-ahead 技术虽然解决了非对称 Lanczos 算法的中断问题, 然而付出的代价也是不容忽视的: 一是使得实现过程变得更加复杂; 二是得到的投影矩阵已经不是三对角矩阵, 致使 QMRES 方法的运算复杂性大为增加. 因此, 就求解线性方程组而言, 通常是放弃使用 Look-ahead, 而是在遇到中断时简单地采用重新开始的办法.

4.5.2 QMRES 方法

从 $v_1 = r_0 / \|r_0\|_2$ 出发, 用算法 4.13 计算得到分解

$$AV_k = V_k T_k + \beta_k v_{k+1} e_k^T = V_{k+1} \tilde{T}_k \quad (4.178)$$

之后, QMRES 方法的下一步是求解最小二乘问题

$$\min \{ \|\beta e_1 - \tilde{T}_k z\|_2 : z \in \mathbb{R}^k \}, \quad (4.179)$$

式中: $\beta = \|r_0\|_2$. 一旦这样的 z_k 求得, 则所寻求的 x_k 就是 $x_k = x_0 + V_k z_k$.

极小化问题 (4.179) 是一个系数矩阵为三对角矩阵的最小二乘问题, 仍然用 QR 分解方法来求解.

由于 \tilde{T}_k 具有式 (4.177) 所示的形状, 故可以计算 k 个 Givens 变换 G_1, G_2, \dots, G_k , 使得

$$G_k G_{k-1} \cdots G_2 G_1 \tilde{T}_k = \begin{bmatrix} R_k \\ O \end{bmatrix}, \quad (4.180)$$

式中:

$$G_i = \text{diag} \left(I_{i-1}, \begin{bmatrix} c_i & s_i \\ -s_i & c_i \end{bmatrix}, I_{k-i} \right) \in \mathbb{R}^{(k+1) \times (k+1)}, \quad c_i^2 + s_i^2 = 1,$$

$$R_k = \begin{bmatrix} \sigma_1 & \delta_1 & \varepsilon_1 & & \\ & \sigma_2 & \delta_2 & \ddots & \\ & & \ddots & \ddots & \varepsilon_{k-2} \\ & & & \sigma_{k-1} & \delta_{k-1} \\ & & & & \sigma_k \end{bmatrix}, \quad (4.181)$$

而且 $\beta_i \neq 0$ 蕴涵着 $\sigma_i \neq 0, i = 1, 2, \dots, k$, 从而 R_k 是非奇异的.

令

$$G = G_k G_{k-1} \cdots G_2 G_1, \quad \begin{bmatrix} t_k \\ \rho_k \end{bmatrix} = G(\beta e_1), \quad t_k = (\tau_1, \tau_2, \dots, \tau_k)^T, \quad (4.182)$$

则 G 是 $k+1$ 阶正交矩阵, 而且直接计算有

$$\begin{aligned} \tau_1 &= \beta c_1, \quad \tau_i = (-1)^{i-1} \beta s_1 s_2 \cdots s_{i-1} c_i, \quad i = 2, 3, \dots, k, \\ \rho_k &= (-1)^k \beta s_1 s_2 \cdots s_k. \end{aligned} \quad (4.183)$$

这样, 利用式 (4.180) 和式 (4.182), 有

$$\begin{aligned} \|\tilde{T}_k z - \beta e_1\|_2^2 &= \|G(\tilde{T}_k z - \beta e_1)\|_2^2 = \left\| \begin{bmatrix} R_k \\ 0 \end{bmatrix} z - \begin{bmatrix} t_k \\ \rho_k \end{bmatrix} \right\|_2^2 \\ &= \|R_k z - t_k\|_2^2 + \rho_k^2 \end{aligned}$$

对任意的 $z \in \mathbb{R}^k$ 成立. 由此立即知道, 最小三乘问题 (4.179) 有唯一解

$$z_k = R_k^{-1} t_k, \quad (4.184)$$

而且有

$$\|\tilde{T}_k z_k - \beta e_1\|_2 = |\rho_k|. \quad (4.185)$$

由式 (4.184) 求得 z_k 之后, 就可算出所求的 x_k 为 $x_k = x_0 + V_k z_k$.

由式 (4.167) 和式 (4.185), 得

$$\|r_k\|_2 = \|b - Ax_k\|_2 \leq \|V_{k+1}\|_2 \|\beta e_1 - \tilde{T}_k z_k\|_2 = \|V_{k+1}\|_2 |\rho_k|,$$

故在实际计算时, 可以用

$$|\rho_k|/\beta \leq \varepsilon \quad (4.186)$$

作为迭代终止的准则, 其中 $\varepsilon > 0$ 是给定的误差要求.

由式 (4.183) 可知, ρ_k 的值并不需要 z_k 和 x_k 的信息. 因此, 只需在 ρ_k 满足式 (4.186) 之后, 再去计算 z_k 和 x_k 即可.

将式 (4.184) 代入 $x_k = x_0 + V_k z_k$, 得

$$x_k = x_0 + V_k R_k^{-1} t_k = x_0 + P_k t_k, \quad (4.187)$$

式中: $P_k = V_k R_k^{-1}$. 这样, 只要将 P_k 算出, 就可以通过式 (4.187) 来计算 x_k . 令 $P_k = [p_1, p_2, \dots, p_k]$, 则比较 $P_k R_k = V_k$ 两边的每一列, 得

$$\begin{aligned} \sigma_1 p_1 &= v_1, \\ \delta_1 p_1 + \sigma_2 p_2 &= v_2, \end{aligned}$$

$$\varepsilon_{i-2}\mathbf{p}_{i-2} + \delta_{i-1}\mathbf{p}_{i-1} + \sigma_i\mathbf{p}_i = \mathbf{v}_i, \quad i = 3, 4, \dots, k.$$

由此可求得 \mathbf{P}_k 的列向量为

$$\begin{cases} \mathbf{p}_1 = \mathbf{v}_1/\sigma_1, \\ \mathbf{p}_2 = (\mathbf{v}_2 - \delta_1\mathbf{p}_1)/\sigma_2, \\ \mathbf{p}_i = (\mathbf{v}_i - \varepsilon_{i-2}\mathbf{p}_{i-2} - \delta_{i-1}\mathbf{p}_{i-1})/\sigma_i, \quad i = 3, 4, \dots, k. \end{cases} \quad (4.188)$$

下面借助式 (4.183), 式 (4.187) 和式 (4.188) 导出计算 \mathbf{x}_k 的递推公式. 首先注意到, 若非对称 Lanczos 分解的长度由 k 增加到 $k+1$, 则有

$$\tilde{\mathbf{T}}_{k+1} = \left[\begin{array}{c|c} \tilde{\mathbf{T}}_k & \tilde{\mathbf{t}}_{k+1} \\ \hline \mathbf{0} & \beta_{k+1} \end{array} \right], \quad \tilde{\mathbf{t}}_{k+1} = (0, \dots, 0, \gamma_k, \alpha_{k+1})^T,$$

于是有

$$\mathbf{R}_{k+1} = \left[\begin{array}{c|c} \mathbf{R}_k & \tilde{\mathbf{r}}_{k+1} \\ \hline \mathbf{0} & \sigma_{k+1} \end{array} \right], \quad \tilde{\mathbf{r}}_{k+1} = (0, \dots, 0, \varepsilon_{k-1}, \delta_k)^T, \quad (4.189)$$

式中:

$$\begin{aligned} \varepsilon_{k-1} &= s_{k-1}\gamma_k, \quad \hat{\gamma}_k = c_{k-1}\gamma_k, \\ \delta_k &= c_k\hat{\gamma}_k + s_k\alpha_{k+1}, \quad \tilde{\alpha}_{k+1} = -s_k\hat{\gamma}_k + c_k\alpha_{k+1}. \end{aligned}$$

上面的四个等式由如下 Givens 变换得到:

$$\begin{aligned} \tilde{\mathbf{T}}_{k+1} &= \left[\begin{array}{cccc|cccc} \alpha_1 & \gamma_1 & & & & & & \\ \beta_1 & \alpha_2 & \gamma_2 & & & & & \\ & \ddots & \ddots & \ddots & & & & \\ & & \beta_{k-3} & \alpha_{k-2} & \gamma_{k-2} & & & \\ & & & \beta_{k-2} & \alpha_{k-1} & \gamma_{k-1} & & \\ & & & & \beta_{k-1} & \alpha_k & \gamma_k & \\ & & & & & \beta_k & \alpha_{k+1} & \\ \hline & & & & & & \beta_{k+1} & \end{array} \right] \xrightarrow{\mathbf{G}_1} \left[\begin{array}{ccc|cccc} \sigma_1 & \delta_1 & \varepsilon_1 & & & & & \\ & \tilde{\alpha}_2 & \hat{\gamma}_2 & & & & & \\ & & \ddots & \ddots & \ddots & & & \\ & & & \alpha_{k-2} & \gamma_{k-2} & & & \\ & & & & \beta_{k-2} & \alpha_{k-1} & \gamma_{k-1} & \\ & & & & & \beta_{k-1} & \alpha_k & \gamma_k \\ & & & & & & \beta_k & \alpha_{k+1} \\ \hline & & & & & & & \beta_{k+1} \end{array} \right] \\ &\xrightarrow{\mathbf{G}_2 \dots \mathbf{G}_{k-2}} \left[\begin{array}{ccc|cccc} \sigma_1 & \delta_1 & \varepsilon_1 & & & & & \\ & \sigma_2 & \delta_2 & \varepsilon_2 & & & & \\ & & \ddots & \ddots & \ddots & & & \\ & & & \sigma_{k-2} & \delta_{k-2} & \varepsilon_{k-2} & & \\ & & & & \tilde{\alpha}_{k-1} & \hat{\gamma}_{k-1} & & \\ & & & & & \beta_{k-1} & \alpha_k & \gamma_k \\ & & & & & & \beta_k & \alpha_{k+1} \\ \hline & & & & & & & \beta_{k+1} \end{array} \right] \end{aligned}$$

$$\xrightarrow{G_{k-1}} \left[\begin{array}{ccc|ccc} \sigma_1 & \delta_1 & \varepsilon_1 & & & \\ & \sigma_2 & \delta_2 & \varepsilon_2 & & \\ & & \ddots & \ddots & \ddots & \\ & & & \sigma_{k-2} & \delta_{k-2} & \varepsilon_{k-2} \\ & & & & \sigma_{k-1} & \delta_{k-1} & \varepsilon_{k-1} \\ & & & & & \tilde{\alpha}_k & \hat{\gamma}_k \\ & & & & & \beta_k & \alpha_{k+1} \\ \hline & & & & & & \beta_{k+1} \end{array} \right] \xrightarrow{G_k} \left[\begin{array}{ccc|ccc} \sigma_1 & \delta_1 & \varepsilon_1 & & & \\ & \sigma_2 & \delta_2 & \varepsilon_2 & & \\ & & \ddots & \ddots & \ddots & \\ & & & \sigma_{k-2} & \delta_{k-2} & \varepsilon_{k-2} \\ & & & & \sigma_{k-1} & \delta_{k-1} & \varepsilon_{k-1} \\ & & & & & \sigma_k & \delta_k \\ & & & & & & \tilde{\alpha}_{k+1} \\ \hline & & & & & & \beta_{k+1} \end{array} \right],$$

即由

$$\begin{bmatrix} c_{k-1} & s_{k-1} \\ -s_{k-1} & c_{k-1} \end{bmatrix} \begin{bmatrix} 0 \\ \gamma_k \end{bmatrix} = \begin{bmatrix} \varepsilon_{k-1} \\ \hat{\gamma}_k \end{bmatrix}, \quad \begin{bmatrix} c_k & s_k \\ -s_k & c_k \end{bmatrix} \begin{bmatrix} \hat{\gamma}_k \\ \alpha_{k+1} \end{bmatrix} = \begin{bmatrix} \delta_k \\ \tilde{\alpha}_{k+1} \end{bmatrix}$$

得到. 而 σ_{k+1} 是在确定第 $k+1$ 个 Givens 变换 G_{k+1} 时得到的, 即计算 $c_{k+1} = \cos \theta_{k+1}$ 和 $s_{k+1} = \sin \theta_{k+1}$, 使得

$$\begin{bmatrix} c_{k+1} & s_{k+1} \\ -s_{k+1} & c_{k+1} \end{bmatrix} \begin{bmatrix} \tilde{\alpha}_{k+1} \\ \beta_{k+1} \end{bmatrix} = \begin{bmatrix} \sigma_{k+1} \\ 0 \end{bmatrix}. \quad (4.190)$$

由式 (4.183) 可知 $\mathbf{t}_{k+1} = (\mathbf{t}_k^T, \tau_{k+1})^T$, 其中

$$\tau_{k+1} = (-1)^k \beta s_1 s_2 \cdots s_k c_{k+1} = \rho_k c_{k+1}, \quad (4.191)$$

而

$$\rho_{k+1} = (-1)^{k+1} \beta s_1 s_2 \cdots s_k s_{k+1} = -\rho_k s_{k+1}. \quad (4.192)$$

由式 (4.188) 和式 (4.189), 有

$$\mathbf{P}_{k+1} = \mathbf{V}_{k+1} \mathbf{R}_{k+1}^{-1} = [\mathbf{V}_k, \mathbf{v}_{k+1}] \left[\begin{array}{c|c} \mathbf{R}_k^{-1} & -\mathbf{R}_k^{-1} \tilde{\mathbf{r}}_{k+1} / \sigma_{k+1} \\ \hline \mathbf{0} & 1 / \sigma_{k+1} \end{array} \right] = [\mathbf{P}_k, \mathbf{p}_{k+1}], \quad (4.193)$$

这里

$$\begin{aligned} \mathbf{p}_{k+1} &= (\mathbf{v}_{k+1} - \mathbf{V}_k \mathbf{R}_k^{-1} \tilde{\mathbf{r}}_{k+1}) / \sigma_{k+1} = (\mathbf{v}_{k+1} - \mathbf{P}_k \tilde{\mathbf{r}}_{k+1}) / \sigma_{k+1} \\ &= (\mathbf{v}_{k+1} - \varepsilon_{k-1} \mathbf{p}_{k-1} - \delta_k \mathbf{p}_k) / \sigma_{k+1}. \end{aligned} \quad (4.194)$$

从而有

$$\mathbf{x}_{k+1} = \mathbf{x}_0 + \mathbf{P}_{k+1} \mathbf{t}_{k+1} = \mathbf{x}_0 + [\mathbf{P}_k, \mathbf{p}_{k+1}] \begin{bmatrix} \mathbf{t}_k \\ \tau_{k+1} \end{bmatrix} = \mathbf{x}_k + \tau_{k+1} \mathbf{p}_{k+1}, \quad (4.195)$$

这就得到了 \mathbf{x}_k 的递推公式.

算法 4.14 (QMRES 方法) 给定 n 阶非奇异的实矩阵 A , n 维向量 b , 初始向量 x_0 和允许误差 $\varepsilon > 0$. 本算法计算向量 x_k , 使得 $\|r_k\|_2/\|r_0\|_2 \leq \varepsilon$, 其中 $r_k = b - Ax_k$.

选取 x_0 ; 计算 $r_0 = b - Ax_0$; $\beta = \|r_0\|_2$; $v_1 = r_0/\beta$;

取一个 $w_1 \in \mathbb{R}^n$, 使得 $w_1^T v_1 = 1$ (例如可取 $w_1 = v_1$);

$\alpha_1 = w_1^T A v_1$; $v = A v_1 - \alpha_1 v_1$; $w = A^T w_1 - \alpha_1 w_1$;

$\omega = w^T v$;

if $\omega = 0$

$x_1 = x_0 + r_0/\alpha_1$; 结束

else

$\beta_1 = \sqrt{|\omega|}$; $\gamma_1 = \omega/\beta_1$;

$v_2 = v/\beta_1$; $w_2 = w/\gamma_1$;

end

确定 $c_1 = \cos \theta_1$ 和 $s_1 = \sin \theta_1$, 使得

$$\begin{bmatrix} c_1 & s_1 \\ -s_1 & c_1 \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \beta_1 \end{bmatrix} = \begin{bmatrix} \sigma_1 \\ 0 \end{bmatrix};$$

$p_1 = v_1/\sigma_1$; $\rho_1 = -\beta s_1$; $\tau_1 = \beta c_1$;

$p_0 = 0$; $c_0 = 1$; $s_0 = 0$; $x_1 = x_0 + \tau_1 p_1$; $k = 1$;

while ($|\rho_k| > \beta\varepsilon$)

$\alpha_{k+1} = w_{k+1}^T A v_{k+1}$;

$v = A v_{k+1} - \alpha_{k+1} v_{k+1} - \gamma_k v_k$;

$w = A^T w_{k+1} - \alpha_{k+1} w_{k+1} - \beta_k w_k$;

$\omega_{k+1} = w^T v$;

if $\omega_{k+1} = 0$

stop

else

$\beta_{k+1} = \sqrt{|\omega_{k+1}|}$; $\gamma_{k+1} = \omega_{k+1}/\beta_{k+1}$;

$v_{k+2} = v/\beta_{k+1}$; $w_{k+2} = w/\gamma_{k+1}$;

end

$\varepsilon_{k-1} = s_{k-1} \gamma_k$; $\hat{\gamma}_k = c_{k-1} \gamma_k$;

$\delta_k = c_k \hat{\gamma}_k + s_k \alpha_{k+1}$; $\tilde{\alpha}_{k+1} = -s_k \hat{\gamma}_k + c_k \alpha_{k+1}$;

确定 $c_{k+1} = \cos \theta_{k+1}$ 和 $s_{k+1} = \sin \theta_{k+1}$, 使得

$$\begin{bmatrix} c_{k+1} & s_{k+1} \\ -s_{k+1} & c_{k+1} \end{bmatrix} \begin{bmatrix} \tilde{\alpha}_{k+1} \\ \beta_{k+1} \end{bmatrix} = \begin{bmatrix} \sigma_{k+1} \\ 0 \end{bmatrix};$$

$\tau_{k+1} = \rho_k c_{k+1}$; $\rho_{k+1} = -\rho_k s_{k+1}$;

$p_{k+1} = (v_{k+1} - \varepsilon_{k-1} p_{k-1} - \delta_k p_k)/\sigma_{k+1}$;

$x_{k+1} = x_k + \tau_{k+1} p_{k+1}$;

$k = k + 1$;

end

注 4.7 (1) 该算法只需存储 10 个 n 维向量即可, 并不随着 k 的增加而增加.

(2) 同样可以考虑预处理 QMRES 方法 (记为 PQMRES). 选定预处理矩阵 M , 然后将 QMRES 方法应用到求解方程组 $M^{-1}Ax = M^{-1}b$. 这只需要在算法 4.14 中将非对称 Lanczos 过程中的矩阵 A 替换成 $M^{-1}A$ 及初始残差 $r_0 = b - Ax_0$ 替换成 $r_0 = M^{-1}(b - Ax_0)$ 即可.

例 4.11 给定线性方程组的系数矩阵 A 和右端项分别为

$$A = \begin{bmatrix} 4 & -1 & & & \\ -2 & 4 & & & \\ & & \ddots & & \\ & & & \ddots & -1 \\ & & & -2 & 4 \end{bmatrix}, \quad b = A \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix},$$

其中 $n = 1000$. 显然, 该方程组的真解为 $x^* = (1, 1, \dots, 1)^T$. 应用 QMRES 方法到该方程组上 ($\varepsilon = 10^{-10}$), 迭代在 25 步后满足终止条件, 计算得到的近似解 $\tilde{x} = x_{25}$ 满足终止条件, 但

$$\|\tilde{x} - x^*\|_2 = 3.0263 \times 10^{-6}, \quad \|b - A\tilde{x}\|_2 = 9.8572 \times 10^{-6}.$$

若选取预处理矩阵为 $M = \text{tril}(A)$ (即 A 的下三角部分). 应用 PQMRES 方法到该方程组上 ($\varepsilon = 10^{-10}$), 迭代在 12 步后满足终止条件, 计算得到的近似解 $\tilde{x} = x_{12}$ 满足终止条件, 但

$$\|\tilde{x} - x^*\|_2 = 3.4207 \times 10^{-5}, \quad \|b - A\tilde{x}\|_2 = 2.3047 \times 10^{-4}.$$

若将 GMRES 法应用于此例, 迭代在 40 步后满足终止条件, 计算得到的近似解 $\tilde{x} = x_{40}$ 满足

$$\|\tilde{x} - x^*\|_2 = 1.5159 \times 10^{-9}, \quad \|b - A\tilde{x}\|_2 = 2.1777 \times 10^{-9}.$$

迭代过程的收敛轨迹如图 4.12 所示, 其中横坐标为迭代步数 k , 纵坐标为相对残差 $\|r_k\|_2 / \|r_0\|_2$, 这里 r_k 是第 k 步得到的残差向量. 值得注意的是, Jacobi 预条件子对本例无效 (矩阵 A 的对角元素相同), Gauss-Seidel 预条件子虽然对迭代的加速很明显, 但在计算时间上没有任何优势.

4.6 LSQR 方法

本节讨论用 LSQR 方法求解 n 阶非对称方程组

$$Ax = b,$$

或当 $A \in \mathbb{R}^{m \times n}$ ($m > n$) 时, 求解最小二乘问题

$$\min \|Ax - b\|_2. \quad (4.196)$$

这类方法基于 Lanczos 双对角化的思想, 是求解非对称方程组尤其是超定方程组 (线性最小二乘问题) 的一类十分有效的子空间方法.

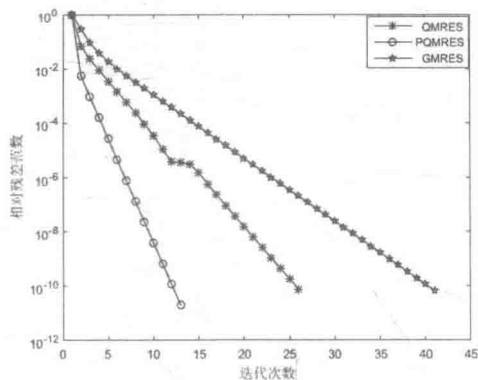


图 4.12 QMRES, PQMRES 和 GMRES 方法的收敛特性

4.6.1 Lanczos 双对角化方法

设 $A \in \mathbb{R}^{m \times n}$ ($m \geq n$). 给定初始向量 $u_1 \in \mathbb{R}^m$ 满足 $\|u_1\|_2 = 1$. 双对角化方法产生 m 维向量组 u_1, u_2, \dots 和 n 维向量组 v_1, v_2, \dots 如下:

$$\alpha_i v_i = A^T u_i - \beta_i v_{i-1} \quad (\beta_1 v_0 = 0), \quad (4.197)$$

$$\beta_{i+1} u_{i+1} = A v_i - \alpha_i u_i, \quad i = 1, 2, \dots, \quad (4.198)$$

式中: $\alpha_i, \beta_{i+1} \geq 0$, 且使 $\|u_{i+1}\|_2 = \|v_i\|_2 = 1$. 因此, 上述双对角方法可如下实现:

算法 4.15 (Lanczos 双对角化方法) 给定矩阵 $A \in \mathbb{R}^{m \times n}$ ($m \geq n$), 初始向量 $u_1 \in \mathbb{R}^m$ 满足 $\|u_1\|_2 = 1$, $v_0 = 0$ 和 $\beta_1 = 0$.

for $i = 1 : k$

$$\hat{v}_i = A^T u_i - \beta_i v_{i-1};$$

$$\alpha_i = \|\hat{v}_i\|_2; \quad v_i = \hat{v}_i / \alpha_i;$$

$$\hat{u}_{i+1} = A v_i - \alpha_i u_i;$$

$$\beta_{i+1} = \|\hat{u}_{i+1}\|_2; \quad u_{i+1} = \hat{u}_{i+1} / \beta_{i+1};$$

end

假设 $\alpha_i, \beta_{i+1} \neq 0, i = 1, 2, \dots, k$, 则上述 Lanczos 双对角化过程可进行到第 k 步, 且若令

$$U_k = [u_1, u_2, \dots, u_k], \quad V_k = [v_1, v_2, \dots, v_k],$$

$$L_k = \begin{bmatrix} \alpha_1 & & & & \\ \beta_2 & \alpha_2 & & & \\ & \ddots & \ddots & & \\ & & \beta_k & \alpha_k & \end{bmatrix}, \quad \tilde{L}_k = \begin{bmatrix} L_k \\ \beta_{k+1} e_k^T \end{bmatrix}, \quad (4.199)$$

则式 (4.197) 和式 (4.198) 可表示为

$$A^T U_k = V_k L_k^T, \quad AV_k = U_k L_k + \beta_{k+1} u_{k+1} e_k^T. \quad (4.200)$$

由此, 得

$$U_k^T AV_k = U_k^T U_k L_k + \beta_{k+1} U_k^T u_{k+1} e_k^T = L_k V_k^T V_k. \quad (4.201)$$

若 $\alpha_{i+1} \neq 0$, 即式 (4.197) 对 $i = k+1$ 成立:

$$\alpha_{k+1} v_{k+1} = A^T u_{k+1} - \beta_{k+1} v_k, \quad (4.202)$$

则由上式和式 (4.200), 得

$$A^T U_{k+1} = V_k \tilde{L}_k^T + \alpha_{k+1} v_{k+1} e_{k+1}^T, \quad AV_k = U_{k+1} \tilde{L}_k. \quad (4.203)$$

由此, 得

$$V_k^T A^T U_{k+1} = V_k^T V_k \tilde{L}_k^T + \alpha_{k+1} V_k^T v_{k+1} e_{k+1}^T = \tilde{L}_k^T U_{k+1}^T U_{k+1}. \quad (4.204)$$

基于以上分析, 可得下面的定理.

定理 4.11 由 Lanczos 双对角化过程产生的两个向量组 $\{u_i\}_{i=1}^k$ 和 $\{v_i\}_{i=1}^k$ 是规范正交的, 即

$$U_k^T U_k = V_k^T V_k = I_k, \quad (4.205)$$

式中: I_k 为 k 阶单位矩阵.

证明 对 k 用归纳法. 当 $k = 1$ 时, 结论显然为真. 设式 (4.205) 对某个 $k \geq 1$ 成立. 利用归纳法假设, 由式 (4.201) 的第 2 个等式可知 $U_k^T u_{k+1} = 0$, 即 u_{k+1} 与 $u_i (i = 1, 2, \dots, k)$ 均正交. 注意到 $\|u_{k+1}\|_2 = 1$, 即得 $U_{k+1}^T U_{k+1} = I_{k+1}$. 再由式 (4.204) 的第 2 个等式可知 $V_k^T v_{k+1} = 0$, 即知 $V_{k+1}^T V_{k+1} = I_{k+1}$. 证毕. \square

由定理 4.11 可知, Lanczos 双对角化过程必在某个 $k \leq \min\{m, n\}$ 中断, 即或在式 (4.200) 中 $\beta_{k+1} u_{k+1} = 0$, 或在式 (4.203) 中 $\alpha_{k+1} v_{k+1} = 0$. 由此可知, 双对角化算法有两个中断状态:

(1) 由 $k \times k$ 阶矩阵 L_k 来表征 (见式 (4.200)):

$$A^T U_k = V_k L_k^T, \quad AV_k = U_k L_k, \quad U_k^T U_k = V_k^T V_k = I_k. \quad (4.206)$$

(2) 由 $(k+1) \times k$ 阶矩阵 \tilde{L}_k 来表征 (见式 (4.203)):

$$A^T U_{k+1} = V_k \tilde{L}_k^T, \quad AV_k = U_{k+1} \tilde{L}_k, \quad U_{k+1}^T U_{k+1} = I_{k+1}, \quad V_k^T V_k = I_k. \quad (4.207)$$

定理 4.12 双对角化过程中断状态是 (1) 的充要条件是 $u_1 \in \mathcal{R}(A)$, 而中断状态是 (2) 的充要条件是 $u_1 \notin \mathcal{R}(A)$.

证明 考虑 \tilde{L}_k 的奇异值分解. 注意到 \tilde{L}_k 是列满秩的, 故有

$$\tilde{L}_k = P_{k+1} \tilde{\Sigma}_k V_k^T, \quad P_{k+1}^T P_{k+1} = I_{k+1}, \quad V_k^T V_k = I_k, \quad (4.208)$$

式中: 矩阵 $\tilde{\Sigma}_k$ 为 $(k+1) \times k$ 阶对角阵,

$$\tilde{\Sigma}_k = \begin{bmatrix} \Sigma_k \\ 0 \end{bmatrix} = \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_k & \\ 0 & \cdots & & 0 \end{bmatrix}, \quad \sigma_i > 0, \quad i = 1, 2, \dots, k.$$

若式 (4.206) 成立, 则结合式 (4.207), 得

$$A^T U_{k+1} P_{k+1} = V_k V_k^T \tilde{\Sigma}_k^T, \quad A V_k V_k = U_{k+1} P_{k+1} \tilde{\Sigma}_k. \quad (4.209)$$

故有

$$A A^T U_{k+1} P_{k+1} = U_{k+1} P_{k+1} \tilde{\Sigma}_k \tilde{\Sigma}_k^T, \quad A^T A V_k V_k = V_k V_k^T \tilde{\Sigma}_k^T \tilde{\Sigma}_k. \quad (4.210)$$

由式 (4.210) 可知, 当式 (4.207) 成立时, A 至少有 k 个奇异值 $\sigma_1, \sigma_2, \dots, \sigma_k$ 非零, 且由式 (4.210) 的第 1 式可知

$$A A^T (U_{k+1} P_{k+1} e_{k+1}) = U_{k+1} P_{k+1} \tilde{\Sigma}_k \tilde{\Sigma}_k^T e_{k+1} = 0,$$

即 A 至少有一个奇异值为零, 且 $U_{k+1} P_{k+1} e_{k+1}$ 是奇异向量. 再由式 (4.209) 的第 1 式可知

$$A^T (U_{k+1} P_{k+1} e_{k+1}) = V_k V_k^T \tilde{\Sigma}_k^T e_{k+1} = 0,$$

因此, $U_{k+1} P_{k+1} e_{k+1} \in \mathcal{N}(A^T)$.

若 $u_1 \in \mathcal{R}(A)$, 则由式 (4.198) 可知 $u_i \in \mathcal{R}(A), \forall i$. 于是有 $U_{k+1} P_{k+1} e_{k+1} \in \mathcal{R}(A)$. 但 $\mathcal{N}(A^T)$ 是 $\mathcal{R}(A)$ 的正交补, 因此, $U_{k+1} P_{k+1} e_{k+1} \notin \mathcal{N}(A^T)$, 这就是说, 式 (4.207) 不可能成立. 所以必成立式 (4.206). 反之, 若 (1) 是中断状态, 即式 (4.206) 成立. 因 L_k 非奇异, 故 $U_k = A V_k L_k^{-1}$, 即 $u_i \in \mathcal{R}(A), i = 1, 2, \dots, k$. 证毕. \square

4.6.2 LSQR 算法

取 $u_1 = r_0 / \beta_1$ ($\beta_1 = \|r_0\|_2$) 作为双对角化算法 4.15 中的初始向量, 用以求解最小二乘问题 (4.196). 令

$$x_k = x_0 + V_k y_k, \quad (4.211)$$

则由式 (4.203), 有

$$\begin{aligned} r_k &= b - A x_k = r_0 - A V_k y_k \\ &= \beta_1 u_1 - U_{k+1} \tilde{L}_k y_k = U_{k+1} (\beta_1 e_1 - \tilde{L}_k y_k). \end{aligned}$$

若令

$$\mathbf{s}_{k+1} = \beta_1 \mathbf{e}_1 - \tilde{\mathbf{L}}_k \mathbf{y}_k, \quad (4.212)$$

则

$$\mathbf{r}_k = \mathbf{U}_{k+1} \mathbf{s}_{k+1}. \quad (4.213)$$

因此极小化问题 (4.196) 等价于

$$\min \|\mathbf{s}_{k+1}\|_2 = \min \|\beta_1 \mathbf{e}_1 - \tilde{\mathbf{L}}_k \mathbf{y}_k\|_2. \quad (4.214)$$

利用正交化方法求解最小二乘问题 (4.214):

作 $(k+1) \times (k+1)$ 阶矩阵 $[\tilde{\mathbf{L}}_k, \beta_1 \mathbf{e}_1]$ 的 QR 分解:

$$\mathbf{G}_k [\tilde{\mathbf{L}}_k, \beta_1 \mathbf{e}_1] = \begin{bmatrix} \mathbf{R}_k & \mathbf{z}_k \\ & \tilde{\zeta}_{k+1} \end{bmatrix} \equiv \begin{bmatrix} \rho_1 & \theta_2 & & & & \zeta_1 \\ & \rho_2 & \theta_3 & & & \zeta_2 \\ & & \ddots & \ddots & & \vdots \\ & & & \rho_{k-1} & \theta_k & \zeta_{k-1} \\ \text{---} & & & & \rho_k & \zeta_k \\ & & & & & \tilde{\zeta}_{k+1} \end{bmatrix}, \quad (4.215)$$

式中: $\mathbf{G}_k \equiv \mathbf{G}_{k,k+1} \cdots \mathbf{G}_{2,3} \mathbf{G}_{1,2}$ 为 Givens 变换的乘积, 用以消去 $\tilde{\mathbf{L}}_k$ 的次对角元 β_2, β_3, \dots .

向量 \mathbf{y}_k 和 \mathbf{s}_{k+1} 可分别由

$$\mathbf{R}_k \mathbf{y}_k = \mathbf{z}_k, \quad \mathbf{z}_k = (\zeta_1, \zeta_2, \dots, \zeta_k)^T \quad (4.216)$$

和

$$\mathbf{s}_{k+1} = \mathbf{G}_k^T \begin{bmatrix} \mathbf{0} \\ \tilde{\zeta}_{k+1} \end{bmatrix} \quad (4.217)$$

求得.

直接求解式 (4.216), \mathbf{y}_k 与 \mathbf{y}_{k-1} 一般并无公共元素, 注意到 $[\mathbf{R}_k, \mathbf{z}_k]$ 可由 $[\mathbf{R}_{k-1}, \mathbf{z}_{k-1}]$ 加入新的一行和一列得到. 由此, 式 (4.211) 与式 (4.216) 有效结合的一个途径是引进矩阵 \mathbf{H}_k :

$$\mathbf{x}_k = \mathbf{x}_0 + \mathbf{V}_k \mathbf{R}_k^{-1} \mathbf{z}_k \equiv \mathbf{x}_0 + \mathbf{H}_k \mathbf{z}_k, \quad (4.218)$$

其中 $\mathbf{H}_k = [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_k]$ 的列可由 $\mathbf{H}_k \mathbf{R}_k = \mathbf{V}_k$ 逐列求得: 令 $\mathbf{h}_0 = \mathbf{0}$, 则成立

$$\theta_k \mathbf{h}_{k-1} + \rho_k \mathbf{h}_k = \mathbf{v}_k,$$

由此可得到递推式

$$\mathbf{h}_k = \frac{1}{\rho_k} (\mathbf{v}_k - \theta_k \mathbf{h}_{k-1}). \quad (4.219)$$

则由式 (4.218), 得

$$\mathbf{x}_k = \mathbf{x}_0 + \sum_{i=1}^{k-1} \zeta_i \mathbf{h}_i + \zeta_k \mathbf{h}_k = \mathbf{x}_{k-1} + \zeta_k \mathbf{h}_k. \quad (4.220)$$

现在需要确定式 (4.215) 的 QR 分解. Givens 变换 $G_{k,k+1}$ 作用于前面 $k-1$ 步变换过的矩阵 $[\tilde{L}_k, \beta_1 e_1]$ 以消去 β_{k+1} , 这给出如下递推关系

$$\begin{bmatrix} c_k & s_k \\ -s_k & c_k \end{bmatrix} \begin{bmatrix} \tilde{\rho}_k & 0 & \tilde{\zeta}_k \\ \beta_{k+1} & \alpha_{k+1} & 0 \end{bmatrix} = \begin{bmatrix} \rho_k & \theta_{k+1} & \zeta_k \\ 0 & \tilde{\rho}_{k+1} & \tilde{\zeta}_{k+1} \end{bmatrix}, \quad (4.221)$$

其中 $\tilde{\rho}_1 \equiv \alpha_1$, $\tilde{\zeta}_1 \equiv \beta_1$. c_k, s_k 是 Givens 变换 $G_{k,k+1}$ 中的元素. $\tilde{\rho}_k$ 和 $\tilde{\zeta}_k$ 是中间变量, 会被 ρ_k 和 ζ_k 所取代. 由 Givens 变换的定义, c_k, s_k 的选取应使得对向量 $[\tilde{\rho}_k, \beta_{k+1}]^T$ 作用后第 2 个分量为零, 且满足 $c_k^2 + s_k^2 = 1$. 为此, 必有

$$c_k = \frac{\tilde{\rho}_k}{\sqrt{\tilde{\rho}_k^2 + \beta_{k+1}^2}}, \quad s_k = \frac{\beta_{k+1}}{\sqrt{\tilde{\rho}_k^2 + \beta_{k+1}^2}}. \quad (4.222)$$

由此, 式 (4.221) 即为

$$\begin{cases} \rho_k = \sqrt{\tilde{\rho}_k^2 + \beta_{k+1}^2}, \quad \theta_{k+1} = s_k \alpha_{k+1}, \quad \zeta_k = c_k \tilde{\zeta}_k, \\ \tilde{\rho}_{k+1} = c_k \alpha_{k+1}, \quad \tilde{\zeta}_{k+1} = -s_k \tilde{\zeta}_k, \quad \tilde{\rho}_1 \equiv \alpha_1, \quad \tilde{\zeta}_1 \equiv \beta_1. \end{cases} \quad (4.223)$$

此外, $G_{k,k+1}$ 在式 (4.221) 中使用完之后即可丢弃而不必存储. 因此由 QR 分解算法产生 R_k, h_k 和 $\tilde{\zeta}_{k+1}$ 的工作量是很小的.

现将 LSQR 算法总结如下, 其中 h_k 用 $z_k = \rho_k h_k$ 代替.

算法 4.16 (LSQR 算法) 给定最小二乘问题 (4.196) 和容许误差限 $\varepsilon > 0$, 取初始向量 x_0 . 本算法计算 x_k , 使得 $\|r_k\|_2 / \|r_0\|_2 \leq \varepsilon$, 其中 $r_k = b - Ax_k$.

选取 x_0 ; 计算 $r_0 = b - Ax_0$; $\beta_1 = \|r_0\|_2$; $u_1 = r_0 / \beta_1$;

$\hat{v}_1 = A^T u_1$; $\alpha_1 = \|\hat{v}_1\|_2$; $v_1 = \hat{v}_1 / \alpha_1$;

$z_1 = v_1$; $\tilde{\zeta}_1 = \beta_1$; $\tilde{\rho}_1 = \alpha_1$; $k = 1$;

while ($\|r_k\|_2 / \|r_0\|_2 > \varepsilon$)

① 双对角化

$$\hat{u}_{k+1} = A v_k - \alpha_k u_k;$$

$$\beta_{k+1} = \|\hat{u}_{k+1}\|_2; \quad u_{k+1} = \hat{u}_{k+1} / \beta_{k+1};$$

$$\hat{v}_{k+1} = A^T u_{k+1} - \beta_{k+1} v_k;$$

$$\alpha_{k+1} = \|\hat{v}_{k+1}\|_2; \quad v_{k+1} = \hat{v}_{k+1} / \alpha_{k+1};$$

② 构造和使用正交变换

$$\rho_k = \sqrt{\tilde{\rho}_k^2 + \beta_{k+1}^2}; \quad c_k = \tilde{\rho}_k / \rho_k; \quad s_k = \beta_{k+1} / \rho_k;$$

$$\theta_{k+1} = s_k \alpha_{k+1}; \quad \tilde{\rho}_{k+1} = c_k \alpha_{k+1};$$

$$\zeta_k = c_k \tilde{\zeta}_k; \quad \tilde{\zeta}_{k+1} = -s_k \tilde{\zeta}_k;$$

③ 更新 x_k , r_k 和 z_k

$$x_{k+1} = x_k + (\zeta_k / \rho_k) z_k;$$

$$r_{k+1} = b - Ax_{k+1} = r_k - (\zeta_k / \rho_k) Az_k;$$

$$z_{k+1} = v_{k+1} - (\theta_{k+1} / \rho_k) z_k;$$

④ 置 $k = k + 1$;

end

注 4.8 由双对角化过程可知

$$v_k \in \text{span}\{A^T u_1, (A^T A) A^T u_1, \dots, (A^T A)^{k-1} A^T u_1\} \equiv \mathcal{K}_k(A^T A, A^T u_1).$$

因此, 由式 (4.211) 和 $u_1 = r_0 / \beta_1$ 可知 LSQR 算法产生的迭代序列 x_k 满足

$$\begin{aligned} x_k &\in x_0 + \text{span}\{A^T r_0, (A^T A) A^T r_0, \dots, (A^T A)^{k-1} A^T r_0\} \\ &\equiv x_0 + \mathcal{K}_k(A^T A, A^T r_0), \end{aligned} \quad (4.224)$$

且使

$$\|r_k\|_2 = \|b - Ax_k\|_2 = \min \{\|b - Ax\|_2 : x \in x_0 + \mathcal{K}_k(A^T A, A^T b)\}. \quad (4.225)$$

LSQR 方法的 MATLAB 程序如下:

```
%LSQR方法程序-mlsqr.m
function [x,k,time,res,resvec]=mlsqr(A,b,x,max_it,tol)
tic; r=b-A*x; mr=norm(r); u=r/mr;
v=A'*u; alpha=norm(v); v=v/alpha;
z=v; zetat=mr; rhot=alpha;
resvec(1)=1; k=0;
while (k<=max_it)
    k=k+1;
    u=A*v-alpha*u; beta=norm(u); u=u/beta;
    v=A'*u-beta*v; alpha=norm(v); v=v/alpha;
    rho=sqrt(rhot^2+beta^2);
    c=rhot/rho; s=beta/rho;
    theta=s*alpha; rhot=c*alpha;
    zeta=c*zetat; zetat=-s*zetat;
    x=x+(zeta/rho)*z; r=b-A*x;
    z=v-(theta/rho)*z;
    res=norm(r)/mr; resvec(k+1)=res;
    if (res<tol), break; end
end
time=toc;
```

例 4.12 假设线性方程组的系数矩阵 A 为

$$A = \begin{bmatrix} 4 & -1 & & & \\ -2 & 4 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -2 & 4 & -1 \\ & & & -2 & 4 \end{bmatrix}.$$

再假设真解 x^* 是分量都为 1 的向量, 从而该方程组的右端项为 $b = Ax^*$. 将 LSQR 算法应用到该线性方程组上, 迭代在 76 步后收敛 ($\varepsilon = 10^{-10}$), 计算得到的近似解 \hat{x} 和真解 x^* 之间的绝对值误差为

$$\|\hat{x} - x^*\|_2 = 1.8330 \times 10^{-9},$$

但计算解 \hat{x} 的残量满足

$$\|b - A\hat{x}\|_2 = 2.5663 \times 10^{-9}.$$

此外, 将 GMRES 方法应用于求解此例, 迭代 40 步满足终止准则. 对于此例, 虽然 LSQR 方法的迭代次数比 GMRES 方法多了将近一倍, 但计算时间不到 GMRES 方法的一半. 迭代过程的收敛轨迹如图 4.13 所示, 其中横坐标为迭代步数 k , 纵坐标为相对残差 $\|r_k\|_2/\|r_0\|_2$, 这里 r_k 是第 k 步得到的残差向量.

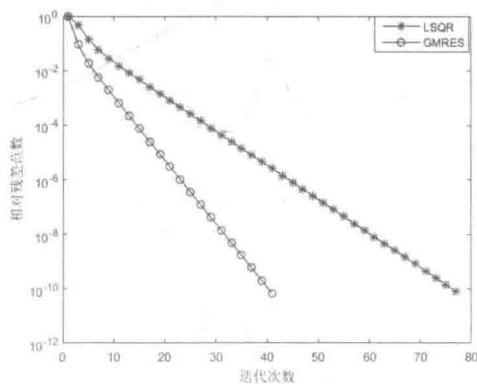


图 4.13 LSQR 算法的收敛特性

4.7 广义共轭残量法

本节考虑求解线性方程组

$$Ax = b, \quad A \in \mathbb{R}^{n \times n}, \quad b \in \mathbb{R}^n \quad (4.226)$$

的广义共轭残量法 (简记为 GCR 方法), 其中系数矩阵 A 是正定的, 即其对称部分 $\frac{1}{2}(A + A^T)$ 是对称正定的.

4.7.1 GCR 方法

类似于对称正定方程组的共轭梯度法, 构造如下的广义共轭残量法.

算法 4.17 (GCR 方法) 给定方程组 (4.226) 和容许误差限 $\varepsilon > 0$, 取初始向量 x_0 . 本算法计算 x_k , 使得 $\|r_k\|_2/\|r_0\|_2 \leq \varepsilon$, 其中 $r_k = b - Ax_k$.

选取 x_0 ; 计算 $r_0 = b - Ax_0$; $p_0 = r_0$; $k = 0$;

while ($\|r_k\|_2/\|r_0\|_2 > \varepsilon$)

$k = k + 1$;

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

$$x_{k+1} = x_k + \alpha_k p_k; r_{k+1} = r_k - \alpha_k Ap_k;$$

for $i = 0, 1, \dots, k$

$$\beta_i^{(k)} = -\frac{(Ar_{k+1}, Ap_i)}{(Ap_i, Ap_i)};$$

end

$$p_{k+1} = r_{k+1} + \sum_{i=0}^k \beta_i^{(k)} p_i;$$

$$Ap_{k+1} = Ar_{k+1} + \sum_{i=0}^k \beta_i^{(k)} Ap_i;$$

end

从算法 4.17 可知, 计算过程中需要存储每一步迭代的 p_i 和 Ap_i ($i = 0, 1, \dots, k$). 其 MATLAB 程序如下:

```
%广义共轭残量法(GCR方法)程序-gcr.m
function [x,k,time,res,resvec]=gcr(A,b,x,max_it,tol)
tic; n=length(b); r=b-A*x; p=r; mr=norm(r);
resvec(1)=1; AP=zeros(n,max_it+1);
P=zeros(n,max_it+1);
AP(:,1)=A*p; P(:,1)=p; k=0;
while (k<max_it)
    k=k+1;
    alpha=(r'*AP(:,k))/(AP(:,k)'*AP(:,k));
    x=x+alpha*P(:,k);
    r=r-alpha*AP(:,k);
    Ar=A*r; s1=zeros(n,1); s2=zeros(n,1);
    for (i=1:k)
        b(i)=-(Ar'*AP(:,i))/(AP(:,i)'*AP(:,i));
        s1=s1+b(i)*P(:,i);
        s2=s2+b(i)*AP(:,i);
    end
end
```

```

end
P(:,k+1)=r+s1;
AP(:,k+1)=Ar+s2;
res=norm(r)/mr; resvec(k+1)=res;
if (res<tol), break; end
end
time=toc;

```

下面列出一些与共轭梯度法 (CG) 相类似的性质:

定理 4.13 设 A 正定, 即 $H = \frac{1}{2}(A + A^T)$ 是对称正定的. 则由算法 4.17 产生的 $\{x_k\}$, $\{r_k\}$ 和 $\{p_k\}$ 具有下列性质:

$$(1) \quad (Ap_k, Ap_l) = 0, \quad k \neq l.$$

$$(2) \quad (r_k, Ap_l) = 0, \quad k > l.$$

$$(3) \quad (r_k, Ap_l) = (r_0, Ap_l), \quad k \leq l.$$

$$(4) \quad (r_k, Ap_k) = (r_k, Ar_k).$$

$$(5) \quad (r_k, Ar_l) = 0, \quad k > l.$$

$$(6) \quad \text{span}\{p_0, p_1, \dots, p_k\} = \text{span}\{p_0, Ap_0, \dots, A^k p_0\} \\ = \text{span}\{r_0, Ar_0, \dots, A^k r_0\} = \text{span}\{r_0, r_1, \dots, r_k\}.$$

$$(7) \quad \text{若 } r_k \neq 0, \text{ 则 } p_k \neq 0.$$

$$(8) \quad x_{k+1} \in \mathcal{K} \equiv x_0 + \text{span}\{p_0, p_1, \dots, p_k\}, \text{ 且使其残量范数 } f(x_{k+1}) \equiv \|r_{k+1}\|_2 = \|b - Ax_{k+1}\|_2 \text{ 在仿射空间 } \mathcal{K} \text{ 中达到极小.}$$

证明 (1) 首先验证 $(Ap_1, Ap_0) = 0$ 成立. 事实上, 由 $\beta_i^{(k)}$ 更新规则, 有

$$(Ap_1, Ap_0) = (A(r_1 + \beta_0^{(0)} p_0), Ap_0) = (Ar_1, Ap_0) + \beta_0^{(0)} (Ap_0, Ap_0) \\ = (Ar_1, Ap_0) - \frac{(Ar_1, Ap_0)}{(Ap_0, Ap_0)} (Ap_0, Ap_0) = 0.$$

利用归纳法, 由 $\beta_i^{(k)}$ 的取法不难验证, 对任意的 $k \neq l$ 都有 $(Ap_k, Ap_l) = 0$.

(2) 对 k 使用归纳法. 当 $k = 0$ 时无需证明. 设 $k \leq s$ 时, 结论为真, 则

$$(r_{s+1}, Ap_l) = (r_s, Ap_l) - \alpha_s (Ap_s, Ap_l).$$

若 $s > l$, 由归纳法假设及结论 (1), 上式右端为 0. 若 $s = l$, 由 α_s 的表达式, 上式右端也为 0.

(3) 对 $k (k \leq l)$ 使用归纳法证明. 当 $k = 0$ 时, 结论平凡地成立. 设 $k = s < l$ 时结论成立, 则由结论 (1) 和归纳法假设, 有

$$(r_{s+1}, Ap_l) = (r_s, Ap_l) - \alpha_s (Ap_s, Ap_l) = (r_s, Ap_l) = (r_0, Ap_l).$$

(4) 利用结论 (2), 有

$$(r_k, Ap_k) = (r_k, Ar_k) + \sum_{i=0}^{k-1} \beta_i^{(k-1)} (r_k, Ap_i) = (r_k, Ar_k).$$

(5) 由算法 4.17, 有

$$r_l = p_l - \sum_{i=0}^{l-1} \beta_i^{(l-1)} p_i.$$

应用结论 (2), 对 $k > l$, 有

$$(r_k, Ar_l) = (r_k, Ap_l) - \sum_{i=0}^{l-1} \beta_i^{(l-1)} (r_k, Ap_i) = 0.$$

(6) 对 k 用归纳法证明. 当 $k = 0$ 时, 结论显然成立. 设对 $k \leq s$ 结论成立, 即

$$r_i, p_i \in \text{span}\{p_0, Ap_0, \dots, A^s p_0\}, \quad i = 0, 1, \dots, s.$$

由算法 4.17, 有

$$p_{s+1} = r_{s+1} + \sum_{i=0}^s \beta_i^{(s)} p_i = r_s - \alpha_s Ap_s + \sum_{i=0}^s \beta_i^{(s)} p_i,$$

故有

$$p_{s+1} \in \text{span}\{p_0, Ap_0, \dots, A^{s+1} p_0\}.$$

由此, 得

$$\text{span}\{p_0, p_1, \dots, p_{s+1}\} \subset \text{span}\{p_0, Ap_0, \dots, A^{s+1} p_0\}.$$

但由结论 (1) 可知 $\{p_i\}$ 线性无关, 故

$$\text{span}\{p_0, p_1, \dots, p_{s+1}\} = \text{span}\{p_0, Ap_0, \dots, A^{s+1} p_0\}.$$

同理, 可证

$$\text{span}\{r_0, r_1, \dots, r_k\} \subset \text{span}\{r_0, Ar_0, \dots, A^k r_0\}.$$

(7) 注意到 $H = \frac{1}{2}(A + A^T)$ 是对称正定的, 若 $r_k \neq 0$, 则由结论 (4), 有

$$(r_k, Ap_k) = (r_k, Ar_k) = (r_k, Hr_k) > 0,$$

故 $(r_k, Ap_k) \neq 0$, 由此 $p_k \neq 0$.

(8) 由算法 4.17, 有

$$x_{k+1} = x_0 + \sum_{i=0}^k \alpha_i p_i \in x_0 + \text{span}\{p_0, p_1, \dots, p_k\}.$$

因此 $f^2(\mathbf{x}_{k+1}) = \|\mathbf{b} - \mathbf{A}\mathbf{x}_{k+1}\|_2^2$ 是 $(\alpha_0, \alpha_1, \dots, \alpha_k)$ 的二次泛函. 应用结论 (1), 得

$$\begin{aligned} f^2(\mathbf{x}_{k+1}) &= \left\| \mathbf{r}_0 - \sum_{i=0}^k \alpha_i \mathbf{A}\mathbf{p}_i \right\|_2^2 \\ &= \|\mathbf{r}_0\|_2^2 - 2 \sum_{i=0}^k \alpha_i (\mathbf{r}_0, \mathbf{A}\mathbf{p}_i) + \sum_{i=0}^k \alpha_i^2 (\mathbf{A}\mathbf{p}_i, \mathbf{A}\mathbf{p}_i). \end{aligned}$$

由 $\frac{\partial f^2(\mathbf{x}_{k+1})}{\partial \alpha_i} = 0$ ($i = 0, 1, \dots, k$), 得

$$\alpha_i = \frac{(\mathbf{r}_0, \mathbf{A}\mathbf{p}_i)}{(\mathbf{A}\mathbf{p}_i, \mathbf{A}\mathbf{p}_i)} = \frac{(\mathbf{r}_i, \mathbf{A}\mathbf{p}_i)}{(\mathbf{A}\mathbf{p}_i, \mathbf{A}\mathbf{p}_i)}, \quad i = 0, 1, \dots, k.$$

这就证明了 $f(\mathbf{x}_{k+1})$ 在仿射空间 \mathcal{K} 中被极小化. 证毕. \square

由定理 4.13 可推得 GCR 方法的有限迭代终止性质.

推论 4.5 设矩阵 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 是正定的 (即其对称部分是对称正定的), 则 GCR 方法至多迭代 n 步就给出精确解.

证明 若当 $k \leq n-1$ 时有 $\mathbf{r}_k = \mathbf{0}$, 这表明 \mathbf{x}_k 已是精确解. 若 $\mathbf{r}_k \neq \mathbf{0}$, $k = 0, 1, \dots, n-1$, 则由定理 4.13 的结论 (7) 可知 $\mathbf{p}_k \neq \mathbf{0}$. 再由结论 (1) 知向量组 $\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_{n-1}$ 线性无关, 因此 $\text{span}\{\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_{n-1}\} = \mathbb{R}^n$. 再由定理 4.13 的结论 (8), \mathbf{x}_n 的残量范数有下述极小化性质:

$$\|\mathbf{r}_n\|_2 = \|\mathbf{b} - \mathbf{A}\mathbf{x}_n\|_2 = \min_{\alpha_i} \left\| \mathbf{r}_0 - \sum_{k=0}^{n-1} \alpha_k \mathbf{A}\mathbf{p}_k \right\|_2.$$

因 \mathbf{A} 非奇异, 故 $\text{span}\{\mathbf{A}\mathbf{p}_0, \mathbf{A}\mathbf{p}_1, \dots, \mathbf{A}\mathbf{p}_{n-1}\} = \mathbb{R}^n$. 从而 $\|\mathbf{r}_n\|_2 = 0$, 故 $\mathbf{x}_n = \mathbf{A}^{-1}\mathbf{b}$. 证毕. \square

现在考虑 GCR 方法作为迭代法时的迭代误差估计和收敛性. 令 \mathcal{P}_k 是次数不超过 k 且满足 $p_k(0) = 1$ 的实系数多项式 $p_k(t)$ 的集合.

定理 4.14 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 正定, 即 $\mathbf{H} = \frac{1}{2}(\mathbf{A} + \mathbf{A}^T)$ 是对称正定的. 则由算法 4.17 产生的残量序列 $\{\mathbf{r}_k\}$ 满足

$$\|\mathbf{r}_k\|_2 \leq \min_{p_k \in \mathcal{P}_k} \|p_k(\mathbf{A})\|_2 \|\mathbf{r}_0\|_2 \leq \left[1 - \frac{\lambda_{\min}^2(\mathbf{H})}{\lambda_{\max}(\mathbf{A}^T \mathbf{A})} \right]^{k/2} \|\mathbf{r}_0\|_2. \quad (4.227)$$

因此, GCR 方法收敛. 若 \mathbf{A} 有完全特征向量集, 则

$$\|\mathbf{r}_k\|_2 \leq \kappa_2(\mathbf{T}) M_k \|\mathbf{r}_0\|_2,$$

式中: \mathbf{T} 为 \mathbf{A} 的 Jordan 标准形的变换矩阵, 而

$$M_k = \min_{p_k \in \mathcal{P}_k} \max_{\lambda \in \sigma(\mathbf{A})} |p_k(\lambda)|.$$

进一步, 若 \mathbf{A} 是规范的, 则

$$\|\mathbf{r}_k\|_2 \leq M_k \|\mathbf{r}_0\|_2.$$

证明 由算法 4.17 和定理 4.13 的结论 (6), GCR 方法产生的残量 \mathbf{r}_k 可表示为 $\mathbf{r}_k = \tilde{p}_k(\mathbf{A})\mathbf{r}_0$, $\tilde{p}_k \in \mathcal{P}_k$. 而任意的 $\mathbf{x} \in \mathbf{x}_0 + \text{span}\{\mathbf{r}_0, \mathbf{A}\mathbf{r}_0, \dots, \mathbf{A}^{k-1}\mathbf{r}_0\}$, 其对应的残量 \mathbf{r} 可表示为 $\mathbf{r} = p_k(\mathbf{A})\mathbf{r}_0$, $p_k \in \mathcal{P}_k$. 于是由定理 4.13 的结论 (8), 有

$$\|\mathbf{r}_k\|_2 = \min_{p_k \in \mathcal{P}_k} \|p_k(\mathbf{A})\mathbf{r}_0\|_2 \leq \min_{p_k \in \mathcal{P}_k} \|p_k(\mathbf{A})\|_2 \|\mathbf{r}_0\|_2,$$

这就证明了式 (4.227) 的第 1 个不等式.

现证明第 2 个不等式. 令

$$p_1(t) = 1 + \alpha t \in \mathcal{P}_1,$$

则 $p_1(t)^k \in \mathcal{P}_k$. 因此

$$\min_{p_k \in \mathcal{P}_k} \|p_k(\mathbf{A})\|_2 \leq \|p_1(\mathbf{A})^k\|_2 \leq \|p_1(\mathbf{A})\|_2^k. \quad (4.228)$$

下面估计 $\|p_1(\mathbf{A})\|_2$. 有

$$\begin{aligned} \|p_1(\mathbf{A})\|_2^2 &= \max_{\mathbf{x} \neq \mathbf{0}} \frac{((\mathbf{I} + \alpha\mathbf{A})\mathbf{x}, (\mathbf{I} + \alpha\mathbf{A})\mathbf{x})}{(\mathbf{x}, \mathbf{x})} \\ &= \max_{\mathbf{x} \neq \mathbf{0}} \left\{ 1 + 2\alpha \frac{(\mathbf{x}, \mathbf{A}\mathbf{x})}{(\mathbf{x}, \mathbf{x})} + \alpha^2 \frac{(\mathbf{A}\mathbf{x}, \mathbf{A}\mathbf{x})}{(\mathbf{x}, \mathbf{x})} \right\}. \end{aligned} \quad (4.229)$$

注意到

$$\frac{(\mathbf{A}\mathbf{x}, \mathbf{A}\mathbf{x})}{(\mathbf{x}, \mathbf{x})} = \frac{(\mathbf{x}, \mathbf{A}^T \mathbf{A}\mathbf{x})}{(\mathbf{x}, \mathbf{x})} \leq \lambda_{\max}(\mathbf{A}^T \mathbf{A}) = \|\mathbf{A}\|_2^2.$$

利用 \mathbf{H} 的对称正定性, 有

$$\frac{(\mathbf{x}, \mathbf{A}\mathbf{x})}{(\mathbf{x}, \mathbf{x})} = \frac{1}{2} \left[\frac{(\mathbf{x}, \mathbf{A}\mathbf{x})}{(\mathbf{x}, \mathbf{x})} + \frac{(\mathbf{x}, \mathbf{A}^T \mathbf{x})}{(\mathbf{x}, \mathbf{x})} \right] = \frac{(\mathbf{x}, \mathbf{H}\mathbf{x})}{(\mathbf{x}, \mathbf{x})} \geq \lambda_{\min}(\mathbf{H}) > 0.$$

若选取 $\alpha < 0$, 则由式 (4.229), 得

$$\|p_1(\mathbf{A})\|_2^2 \leq 1 + 2\lambda_{\min}(\mathbf{H})\alpha + \lambda_{\max}(\mathbf{A}^T \mathbf{A})\alpha^2. \quad (4.230)$$

上式右端当

$$\alpha = -\frac{\lambda_{\min}(\mathbf{H})}{\lambda_{\max}(\mathbf{A}^T \mathbf{A})}$$

时达到极小. 将这个 α 值代入式 (4.230), 得

$$\|p_1(\mathbf{A})\|_2^2 \leq 1 - \frac{\lambda_{\min}^2(\mathbf{H})}{\lambda_{\max}(\mathbf{A}^T \mathbf{A})}.$$

由上式及式 (4.228) 即得式 (4.227) 的第 2 个不等式.

若 A 有完全特征向量集, 即 A 可对角化, 则 $A = T^{-1}AT$ 是对角矩阵, 则

$$\begin{aligned} \|r_k\|_2 &= \min_{p_k \in \mathcal{P}_k} \|Tp_k(A)T^{-1}r_0\|_2 \leq \|T\|_2 \|T^{-1}\|_2 \min_{p_k \in \mathcal{P}_k} \|p_k(A)\|_2 \|r_0\|_2 \\ &= \kappa_2(T) \min_{p_k \in \mathcal{P}_k} \max_{\lambda \in \sigma(A)} |p_k(\lambda)| \cdot \|r_0\|_2. \end{aligned}$$

进一步, 若 A 还是规范的, 则可使 $\kappa_2(T) = 1$. 证毕. □

例 4.13 考虑正定方程组 $Ax = b$, 其中

$$A = \begin{bmatrix} 12 & 3 & 2 & 1 & & & & \\ -3 & 12 & 3 & 2 & 1 & & & \\ -2 & -3 & \ddots & \ddots & \ddots & \ddots & & \\ -1 & -2 & \ddots & \ddots & \ddots & \ddots & 1 & \\ & \ddots & \ddots & \ddots & \ddots & 3 & 2 & \\ & & -1 & -2 & -3 & 12 & 3 & \\ & & & -1 & -2 & -3 & 12 & \end{bmatrix} \in \mathbb{R}^{n \times n}, \quad b = A \begin{bmatrix} 1 \\ 1 \\ \vdots \\ \vdots \\ \vdots \\ 1 \\ 1 \end{bmatrix} \in \mathbb{R}^n.$$

取 $n = 1000$, 将 GCR 方法应用到该线性方程组上, 迭代在 20 步后收敛 ($\varepsilon = 10^{-10}$), 计算得到的近似解 \hat{x} 和真解 x^* 之间的绝对值误差为

$$\|\hat{x} - x^*\|_2 = 2.0725 \times 10^{-9},$$

但计算解 \hat{x} 的残量满足

$$\|b - A\hat{x}\|_2 = 2.7746 \times 10^{-8}.$$

迭代过程的收敛轨迹如图 4.14 所示, 其中横坐标为迭代步数 k , 纵坐标为 $\lg \|r_k\|_2$, 其

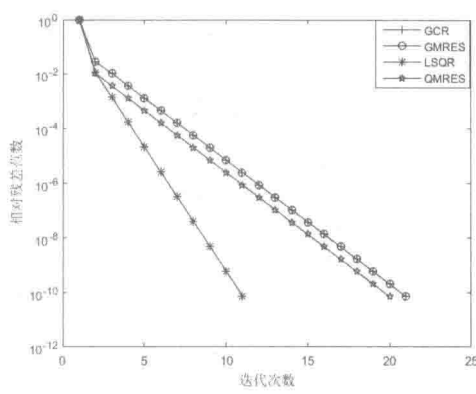


图 4.14 GCR 方法的收敛特性

中 r_k 是第 k 步得到的残量. 此外, 还给出了 GCR 方法与 GMRES, LSQR 及 QMRES 方法的比较, 结果如表 4.3 所示.

从表中可以看出, 对于此例, GCR 算法不如 LSQR 算法有效.

表 4.3 GCR 方法的数值结果

算法	迭代次数	CPU时间	相对残差	绝对误差
GCR	20	0.0028	7.3094e-11	2.0725e-09
GMRES	20	0.0018	7.3094e-11	2.0725e-09
LSQR	10	0.0007	7.7475e-11	2.1967e-09
QMRES	19	0.0010	7.3094e-11	2.0725e-09

4.7.2 GCR(m) 方法

注意到算法 4.17 (GCR 方法) 的执行需要存储所有的方向向量 p_k . 如果在 GCR 方法中周期地重新开始: 每 $m+1$ 次迭代, 用当前的第 i 次重开始的结果 $x_{i(m+1)}$ 作为新的 $i+1$ 次重开始的初始向量, 这时只要存储 m 个方向向量. 这个重新开始 CGR 方法称为 GCR(m). 下面给出重开始 GCR 方法的详细算法步骤.

算法 4.18 (GCR(m) 方法) 给定方程组 (4.226) 和容许误差限 $\varepsilon > 0$, 取初始向量 x_0 及重开始数 m . 本算法计算 x_k , 使得 $\|r_k\|_2/\|r_0\|_2 \leq \varepsilon$, 其中 $r_k = b - Ax_k$.

选取 x_0 ; 计算 $r_0 = b - Ax_0$; $x_1 = x_0$; $r_1 = r_0$; $p_1 = r_0$; $k = 1$;

while ($\|r_k\|_2/\|r_0\| > \varepsilon$)

for $i = 1 : m$

$$\alpha_i = \frac{(r_i, Ap_i)}{(Ap_i, Ap_i)};$$

$$x_{i+1} = x_i + \alpha_i p_i; \quad r_{i+1} = r_i - \alpha_i Ap_i;$$

for $s = 1 : i$

$$\beta_s^{(i)} = -\frac{(Ar_{i+1}, Ap_s)}{(Ap_s, Ap_s)};$$

end

$$p_{i+1} = r_{i+1} + \sum_{s=0}^i \beta_s^{(i)} p_s;$$

$$Ap_{i+1} = Ar_{i+1} + \sum_{s=0}^i \beta_s^{(i)} Ap_s;$$

end

$$x_1 = x_{m+1}; \quad r_1 = r_{m+1}; \quad p_1 = p_{m+1}; \quad Ap_1 = Ap_{m+1};$$

$$k = k + 1;$$

end

GCR(m) 方法的 MATLAB 程序如下:

```
%GCR(m)方法程序-gcrm.m
function [x,out,int,time,res,resvec]=gcrm(A,b,x,restrt,max_it,tol)
tic; n=length(b); m=restrt;
r=b-A*x; p=r; mr=norm(r); k=0;
AP=zeros(n,m+1); P=zeros(n,m+1);
```

```

AP(:,1)=A*p; P(:,1)=p; resvec(1)=1;
while (k<max_it)
    k=k+1;
    for i=1:m
        alpha=(r'*AP(:,i))/(AP(:,i)'*AP(:,i));
        x=x+alpha*P(:,i);
        r=r-alpha*AP(:,i);
        res=norm(r)/mr;
        resvec((k-1)*m+1+i)=res;
        if (res<tol), break; end
        Ar=A*r; s1=zeros(n,1); s2=zeros(n,1);
        for (s=1:i)
            b(s)=-(Ar'*AP(:,s))/(AP(:,s)'*AP(:,s));
            s1=s1+b(s)*P(:,s);
            s2=s2+b(s)*AP(:,s);
        end
        P(:,i+1)=r+s1;
        AP(:,i+1)=Ar+s2;
    end
    if (res<tol),
        out=k; int=i; break;
    end
    AP(:,1)=AP(:,m+1); P(:,1)=P(:,m+1);
end
time=toc;

```

例 4.14 仍考虑例 4.13 中的 A 和 b . 将重开始 GCR 方法应用到该线性方程组上, 取重开始数 $m = 6$, 在外迭代第 4 步内迭代第 2 步后收敛 ($\varepsilon = 10^{-10}$), 计算得到的近似解 \hat{x} 和真解 x^* 之间的绝对值误差为

$$\|\hat{x} - x^*\|_2 = 2.0725 \times 10^{-9},$$

但计算解 \hat{x} 的残量满足

$$\|b - A\hat{x}\|_2 = 2.7746 \times 10^{-8}.$$

迭代过程的收敛轨迹如图 4.15 所示, 其中横坐标为迭代步数 k , 纵坐标为相对残差 $\|r_k\|_2/\|r_0\|_2$, 这里 r_k 是第 k 步得到的残差向量.

若按公式

$$\text{总迭代次数} = (\text{外迭代次数} - 1) \times \text{重开始数} + \text{内迭代数}$$

计算, 本例的总迭代次数为 20, 这跟例 4.13 中 GCR 方法一致.

此外, 还比较了 GCR(m) 与 GMRES(m) 算法. 从图 4.15 可以看出, 对于此例, GCR(m) 算法跟 GMRES(m) 算法的数值效果相仿.

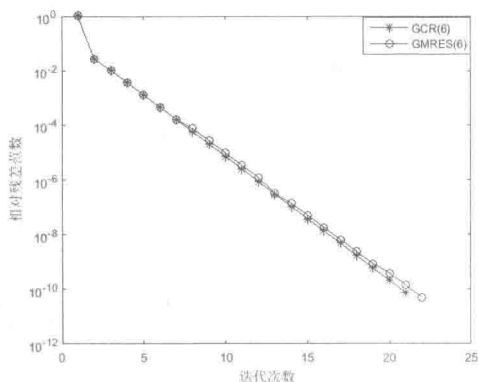


图 4.15 GCR(m) 方法的收敛特性

关于 GCR(m) 方法的收敛性, 有下面的结论.

定理 4.15 在定理 4.14 的假设下, 若 $\{\mathbf{r}_k\}$ 是 GCR(m) 的残量序列, 则

$$\|\mathbf{r}_{i(m+1)}\|_2 \leq \left[\min_{q_{m+1} \in \mathcal{P}_{m+1}} \|q_{m+1}(\mathbf{A})\|_2 \right]^i \|\mathbf{r}_0\|_2, \quad (4.231)$$

由此, 得

$$\|\mathbf{r}_k\|_2 \leq \left[1 - \frac{\lambda_{\min}^2(\mathbf{H})}{\lambda_{\max}(\mathbf{A}^\top \mathbf{A})} \right]^{k/2} \|\mathbf{r}_0\|_2. \quad (4.232)$$

因此, GCR(m) 方法收敛.

证明 由定理 4.14, 有

$$\begin{aligned} \|\mathbf{r}_{i(m+1)}\|_2 &\leq \left[\min_{q_{m+1} \in \mathcal{P}_{m+1}} \|q_{m+1}(\mathbf{A})\|_2 \right] \|\mathbf{r}_{(i-1)(m+1)}\|_2 \\ &\leq \cdots \leq \left[\min_{q_{m+1} \in \mathcal{P}_{m+1}} \|q_{m+1}(\mathbf{A})\|_2 \right]^i \|\mathbf{r}_0\|_2. \end{aligned}$$

令 $k = i(m+1) + s$, $0 \leq s \leq m$, 再由定理 4.14, 有

$$\|\mathbf{r}_k\|_2 \equiv \|\mathbf{r}_{i(m+1)+s}\|_2 \leq \left[1 - \frac{\lambda_{\min}^2(\mathbf{H})}{\lambda_{\max}(\mathbf{A}^\top \mathbf{A})} \right]^{s/2} \|\mathbf{r}_{i(m+1)}\|_2. \quad (4.233)$$

由式 (4.231) 和定理 4.14, 得

$$\|\mathbf{r}_{i(m+1)}\|_2 \leq \left[1 - \frac{\lambda_{\min}(\mathbf{H})^2}{\lambda_{\max}(\mathbf{A}^\top \mathbf{A})} \right]^{i(m+1)/2} \|\mathbf{r}_0\|_2.$$

将上式代入式 (4.233) 即得式 (4.232). 证毕. \square

4.8 投影类方法

前几节所介绍的方法都是残量极小化类型的方法, 本节介绍第二类方法, 即残量正交化方法. 考虑如下的线性方程组

$$Ax = b, \quad (4.234)$$

式中: 矩阵 $A \in \mathbb{R}^{n \times n}$ 和向量 $b \in \mathbb{R}^n$ 是已经给定的, 而 $x \in \mathbb{R}^n$ 是待求的未知向量. 残量正交化方法是求 $x_k \in x_0 + \mathcal{K}_k(A, r_0)$, 使得

$$r_k \perp \mathcal{X}_k, \quad (4.235)$$

其中 $r_k = b - Ax_k$, 而 $\mathcal{X}_k \subset \mathbb{R}^n$ 是一个适当选择的 k 维子空间. 通常称式 (4.235) 为 Galerkin 条件.

设

$$\begin{aligned} \mathcal{X}_k &= \text{span}\{w_1, w_2, \dots, w_k\}, \\ \mathcal{K}_k(A, r_0) &= \text{span}\{v_1, v_2, \dots, v_k\}, \end{aligned}$$

并且记

$$V_k = [v_1, v_2, \dots, v_k], \quad W_k = [w_1, w_2, \dots, w_k],$$

则容易导出, 求 $x_k \in x_0 + \mathcal{K}_k(A, r_0)$ 满足式 (4.235) 就等价于求解线性方程组

$$H_k z = f_k, \quad (4.236)$$

式中:

$$H_k = W_k^T A V_k, \quad f_k = W_k^T r_0. \quad (4.237)$$

从几何上来看, 当 $W_k^T V_k = I_k$ 时, H_k 正好是 A 沿 \mathcal{X}_k^\perp 到 $\mathcal{K}_k(A, r_0)$ 上的投影 (见注 4.6), 故称式 (4.236) 为投影方程, 而称这类方法为投影类方法.

通常 $k \ll n$, 故式 (4.236) 有多种方法来求解. 这样, 实现这一方法的关键是如何选择子空间 \mathcal{X}_k 以及如何有效地计算 $\mathcal{K}_k(A, r_0)$ 和 \mathcal{X}_k 的基向量.

本节取 $\mathcal{X}_k = \mathcal{K}_k(A^T, \tilde{r}_0)$, $\tilde{r}_0^T r_0 \neq 0$, 并给出这方面的三类典型的方法: 双共轭梯度法 (BCG 方法)、共轭梯度平方法 (CGS 方法) 和稳定化共轭梯度法 (BCGSTAB 方法). BCG 方法是基础, 另外两个方法是通过改进 BCG 方法而得到的.

4.8.1 BCG 方法

BCG 方法 (双共轭梯度法, Bi-Conjugate Gradient) 是在式 (4.235) 中取子空间 $\mathcal{X}_k = \mathcal{K}_k(A^T, \tilde{r}_0)$ 而得到的一类投影方法. 当然, 可以利用非对称 Lanczos 方法来计算 $\mathcal{K}_k(A, r_0)$ 和 $\mathcal{K}_k(A^T, \tilde{r}_0)$ 的基, 然后以类似于 CG 方法的推导过程导出 BCG 的基本迭代格式. 为了更清楚地展示这类方法与共轭梯度法的联系, 这里采用 CG 方法的一种等价描述来导出 BCG 的迭代格式.

1. CG 方法的等价描述

为了便于参照, 先将求对称正定线性方程组的 CG 方法的迭代格式重新叙述如下:

- (1) $\mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{x}_0$; $\mathbf{p}_{-1} = \mathbf{0}$; $\rho_{-1} = 1$.
- (2) $\rho_k = \mathbf{r}_k^T \mathbf{r}_k$; $\beta_k = \rho_k / \rho_{k-1}$; $\mathbf{p}_k = \mathbf{r}_k + \beta_k \mathbf{p}_{k-1}$;
 $\sigma_k = \mathbf{p}_k^T \mathbf{A} \mathbf{p}_k$; $\alpha_k = \rho_k / \sigma_k$; $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$;
 $\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k$; $k = 0, 1, 2, \dots$.

若迭代进行到了 k 步, 则将产生三组向量:

- ① 近似解向量组 $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_k$;
- ② 残量组 $\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_k$;
- ③ 方向向量组 $\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_k$.

它们具有如下的基本性质:

- ① $\mathbf{r}_i^T \mathbf{r}_j = \rho_i \delta_{ij}$ (残量相互正交);
- ② $\mathbf{p}_i^T \mathbf{A} \mathbf{p}_j = \sigma_i \delta_{ij}$ (方向向量相互 \mathbf{A} -共轭正交);
- ③ $\mathcal{K}_{k+1}(\mathbf{A}, \mathbf{r}_0) = \text{span}\{\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_k\} = \text{span}\{\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_k\}$.

由性质 ③ 可知, 必存在 $\varphi_k, \psi_k \in \mathcal{P}_k$, 使得

$$\mathbf{r}_k = \varphi_k(\mathbf{A})\mathbf{r}_0, \quad \mathbf{p}_k = \psi_k(\mathbf{A})\mathbf{r}_0.$$

若记 $\varphi_0(t) = 1, \psi_{-1}(t) = 0, \vartheta(t) = t$, 则由前面所述的 CG 方法的迭代格式可导出

$$\begin{cases} \varphi_{k+1} = \varphi_k - \alpha_k \vartheta \psi_k, \\ \psi_k = \varphi_k + \beta_k \psi_{k-1}, \quad k = 0, 1, 2, \dots \end{cases} \quad (4.238)$$

现在在线性空间 \mathcal{P}_n 上定义双线性泛函

$$[\varphi, \psi] = \mathbf{r}_0^T \varphi(\mathbf{A}) \psi(\mathbf{A}) \mathbf{r}_0, \quad \forall \varphi, \psi \in \mathcal{P}_n. \quad (4.239)$$

注意, 这样定义的二元泛函除不能由 $[\varphi, \varphi] = 0$ 导出 $\varphi = 0$ 之外, 它具有内积的其他所有性质.

利用式 (4.239) 容易导出 CG 方法迭代格式中的几个参数的表达式, 即

$$\rho_k = [\varphi_k, \varphi_k], \quad \sigma_k = [\psi_k, \vartheta \psi_k]. \quad (4.240)$$

将式 (4.238) 和式 (4.240) 相结合就得到了 CG 方法用多项式语言给出的等价描述:

$$(1) \quad \varphi_0 = 1; \quad \psi_{-1} = 0; \quad \rho_{-1} = 1; \quad (4.241a)$$

$$(2) \quad \rho_k = [\varphi_k, \varphi_k]; \quad \beta_k = \rho_k / \rho_{k-1}; \quad (4.241b)$$

$$\psi_k = \varphi_k + \beta_k \psi_{k-1}; \quad (4.241c)$$

$$\sigma_k = [\psi_k, \vartheta \psi_k]; \quad \alpha_k = \rho_k / \sigma_k; \quad (4.241d)$$

$$\varphi_{k+1} = \varphi_k - \alpha_k \vartheta \psi_k; \quad k = 0, 1, 2, \dots \quad (4.241e)$$

显然, 由上面的迭代格式产生的多项式系列满足

$$[\varphi_i, \varphi_j] = \rho_i \delta_{ij}, \quad [\psi_i, \vartheta \psi_j] = \sigma_i \delta_{ij}. \quad (4.242)$$

2. BCG 方法

类比于式 (4.239), 对给定的非奇异矩阵 \mathbf{A} 以及向量 \mathbf{r}_0 和 $\tilde{\mathbf{r}}_0$, 在 \mathcal{P}_n 上定义双线性泛函

$$\langle \varphi, \psi \rangle = \tilde{\mathbf{r}}_0^T \varphi(\mathbf{A}) \psi(\mathbf{A}) \mathbf{r}_0, \quad \forall \varphi, \psi \in \mathcal{P}_n, \quad (4.243)$$

然后将式 (4.239) 中定义的双线性泛函 $[\cdot, \cdot]$ 换为现在定义的 $\langle \cdot, \cdot \rangle$, 再逐字照搬过来, 便有

$$(1) \quad \varphi_0 = 1; \quad \psi_{-1} = 0; \quad \rho_{-1} = 1; \quad (4.244a)$$

$$(2) \quad \rho_k = \langle \varphi_k, \varphi_k \rangle; \quad \beta_k = \rho_k / \rho_{k-1}; \quad (4.244b)$$

$$\psi_k = \varphi_k + \beta_k \psi_{k-1}; \quad (4.244c)$$

$$\sigma_k = \langle \psi_k, \vartheta \psi_k \rangle; \quad \alpha_k = \rho_k / \sigma_k; \quad (4.244d)$$

$$\varphi_{k+1} = \varphi_k - \alpha_k \vartheta \psi_k; \quad k = 0, 1, 2, \dots. \quad (4.244e)$$

当然, 也可用数学归纳法证明这样迭代产生的 $\{\varphi_i\}_{i=0}^k$ 和 $\{\psi_i\}_{i=0}^k$ 具有性质

$$\langle \varphi_i, \varphi_j \rangle = \rho_i \delta_{ij}, \quad \langle \psi_i, \vartheta \psi_j \rangle = \sigma_i \delta_{ij}. \quad (4.245)$$

现在定义

$$\mathbf{r}_k = \varphi_k(\mathbf{A}) \mathbf{r}_0, \quad \tilde{\mathbf{r}}_k = \varphi_k(\mathbf{A}^T) \tilde{\mathbf{r}}_0, \quad (4.246)$$

$$\mathbf{q}_k = \psi_k(\mathbf{A}) \mathbf{r}_0, \quad \tilde{\mathbf{q}}_k = \psi_k(\mathbf{A}^T) \tilde{\mathbf{r}}_0. \quad (4.247)$$

将式 (4.244) 改成向量形式, 便有

$$(1) \quad \mathbf{r}_0 = \mathbf{b} - \mathbf{A} \mathbf{x}_0; \quad \tilde{\mathbf{r}}_0 \text{ 满足 } \tilde{\mathbf{r}}_0^T \mathbf{r}_0 \neq 0;$$

$$\mathbf{q}_{-1} = \tilde{\mathbf{q}}_{-1} = \mathbf{0}; \quad \rho_{-1} = 1;$$

$$(2) \quad \rho_k = \tilde{\mathbf{r}}_k^T \mathbf{r}_k; \quad \beta_k = \rho_k / \rho_{k-1}; \quad (4.248a)$$

$$\mathbf{q}_k = \mathbf{r}_k + \beta_k \mathbf{q}_{k-1}; \quad (4.248b)$$

$$\tilde{\mathbf{q}}_k = \tilde{\mathbf{r}}_k + \beta_k \tilde{\mathbf{q}}_{k-1}; \quad (4.248c)$$

$$\sigma_k = \tilde{\mathbf{q}}_k^T \mathbf{A} \mathbf{q}_k; \quad \alpha_k = \rho_k / \sigma_k; \quad (4.248d)$$

$$\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{q}_k; \quad (4.248e)$$

$$\tilde{\mathbf{r}}_{k+1} = \tilde{\mathbf{r}}_k - \alpha_k \mathbf{A}^T \tilde{\mathbf{q}}_k; \quad k = 0, 1, 2, \dots. \quad (4.248f)$$

利用数学归纳法可证如下的结果.

定理 4.16 由式 (4.248) 产生的序列 $\{\mathbf{q}_k\}$, $\{\tilde{\mathbf{q}}_k\}$, $\{\mathbf{r}_k\}$ 和 $\{\tilde{\mathbf{r}}_k\}$ 具有如下性质:

$$(1) \quad \tilde{\mathbf{r}}_k^T \mathbf{r}_l = \mathbf{r}_k^T \tilde{\mathbf{r}}_l = \rho_k \delta_{kl}. \quad (4.249a)$$

$$(2) \quad \tilde{\mathbf{q}}_k^T \mathbf{A} \mathbf{q}_l = \mathbf{q}_k^T \mathbf{A} \tilde{\mathbf{q}}_l = \sigma_k \delta_{kl}. \quad (4.249b)$$

$$(3) \quad \mathcal{K}_k(\mathbf{A}, \mathbf{r}_0) = \text{span}\{\mathbf{q}_0, \mathbf{q}_1, \dots, \mathbf{q}_{k-1}\} = \text{span}\{\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_{k-1}\}. \quad (4.249c)$$

$$(4) \quad \mathcal{K}_k(\mathbf{A}^T, \tilde{\mathbf{r}}_0) = \text{span}\{\tilde{\mathbf{q}}_0, \tilde{\mathbf{q}}_1, \dots, \tilde{\mathbf{q}}_{k-1}\} = \text{span}\{\tilde{\mathbf{r}}_0, \tilde{\mathbf{r}}_1, \dots, \tilde{\mathbf{r}}_{k-1}\}. \quad (4.249d)$$

证明 首先对 k 用归纳法证明式 (4.249a) 和式 (4.249b). 显然, 当 $k=0$ 时, 结论是平凡的. 故设 $k+1 \geq 1$. 假设结论对 k 成立. 由式 (4.248c) 和式 (4.248e), 得

$$\begin{aligned} \mathbf{r}_{k+1}^T \tilde{\mathbf{r}}_l &= \mathbf{r}_k^T \tilde{\mathbf{r}}_l - \alpha_k \mathbf{q}_k^T \mathbf{A}^T \tilde{\mathbf{r}}_l \\ &= \mathbf{r}_k^T \tilde{\mathbf{r}}_l - \alpha_k \mathbf{q}_k^T \mathbf{A}^T (\tilde{\mathbf{q}}_l - \beta_l \tilde{\mathbf{q}}_{l-1}). \end{aligned} \quad (4.250)$$

在式 (4.250) 中, 当 $l=k$ 时, 由 α_k 的表达式和归纳假设, 得

$$\begin{aligned} \mathbf{r}_{k+1}^T \tilde{\mathbf{r}}_k &= \mathbf{r}_k^T \tilde{\mathbf{r}}_k - \frac{\tilde{\mathbf{r}}_k^T \mathbf{r}_k}{\tilde{\mathbf{q}}_k^T \mathbf{A} \mathbf{q}_k} (\mathbf{q}_k^T \mathbf{A}^T \tilde{\mathbf{q}}_k - \beta_k \mathbf{q}_k^T \mathbf{A}^T \tilde{\mathbf{q}}_{k-1}) \\ &= \mathbf{r}_k^T \tilde{\mathbf{r}}_k - \frac{\mathbf{r}_k^T \tilde{\mathbf{r}}_k}{\tilde{\mathbf{q}}_k^T \mathbf{A} \mathbf{q}_k} (\tilde{\mathbf{q}}_k^T \mathbf{A} \mathbf{q}_k - \beta_k \tilde{\mathbf{q}}_{k-1}^T \mathbf{A} \mathbf{q}_k) = 0. \end{aligned}$$

当 $l < k$ 时, 由归纳法假设可得 $\mathbf{r}_{k+1}^T \tilde{\mathbf{r}}_l = 0$. 同理, 可证 $\tilde{\mathbf{r}}_{k+1}^T \mathbf{r}_l = 0$ ($l \leq k$). 这就证明了式 (4.249a).

由式 (4.248b) 和式 (4.248f), 得

$$\begin{aligned} \mathbf{q}_{k+1}^T \mathbf{A}^T \tilde{\mathbf{q}}_l &= \mathbf{r}_{k+1}^T \mathbf{A}^T \tilde{\mathbf{q}}_l + \beta_{k+1} \mathbf{q}_k^T \mathbf{A}^T \tilde{\mathbf{q}}_l \\ &= \mathbf{r}_{k+1}^T \left(\frac{\tilde{\mathbf{r}}_l - \tilde{\mathbf{r}}_{l+1}}{\alpha_k} \right) + \beta_{k+1} \mathbf{q}_k^T \mathbf{A}^T \tilde{\mathbf{q}}_l. \end{aligned} \quad (4.251)$$

在式 (4.251) 中, 当 $l=k$ 时, 由 α_k 和 β_{k+1} 的表达式及归纳假设, 得

$$\begin{aligned} \mathbf{q}_{k+1}^T \mathbf{A}^T \tilde{\mathbf{q}}_k &= \mathbf{r}_{k+1}^T \left(\frac{\tilde{\mathbf{r}}_k - \tilde{\mathbf{r}}_{k+1}}{\alpha_k} \right) + \beta_{k+1} \frac{\mathbf{r}_k^T \tilde{\mathbf{r}}_k}{\alpha_k} \\ &= \mathbf{r}_{k+1}^T \left(\frac{\tilde{\mathbf{r}}_k - \tilde{\mathbf{r}}_{k+1}}{\alpha_k} \right) + \frac{\mathbf{r}_{k+1}^T \tilde{\mathbf{r}}_{k+1}}{\mathbf{r}_k^T \tilde{\mathbf{r}}_k} \frac{\mathbf{r}_k^T \tilde{\mathbf{r}}_k}{\alpha_k} \\ &= \frac{\mathbf{r}_{k+1}^T \tilde{\mathbf{r}}_k}{\alpha_k} = 0. \end{aligned}$$

当 $l < k$ 时, 利用式 (4.249a) 和归纳法假设可得 $\mathbf{q}_{k+1}^T \mathbf{A}^T \tilde{\mathbf{q}}_l = 0$. 同理, 可证 $\tilde{\mathbf{q}}_{k+1}^T \mathbf{A} \mathbf{q}_l = 0$ ($l \leq k$). 这就证明了式 (4.249b).

由于

$$\mathcal{K}_k(\mathbf{A}, \mathbf{r}_0) = \text{span}\{\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_{k-1}\}, \quad \mathcal{K}_k(\mathbf{A}^T, \tilde{\mathbf{r}}_0) = \text{span}\{\tilde{\mathbf{r}}_0, \tilde{\mathbf{r}}_1, \dots, \tilde{\mathbf{r}}_{k-1}\},$$

且向量组 $\{\mathbf{q}_0, \mathbf{q}_1, \dots, \mathbf{q}_{k-1}\}$ 与 $\{\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_{k-1}\}$ 及向量组 $\{\tilde{\mathbf{q}}_0, \tilde{\mathbf{q}}_1, \dots, \tilde{\mathbf{q}}_{k-1}\}$ 与 $\{\tilde{\mathbf{r}}_0, \tilde{\mathbf{r}}_1, \dots, \tilde{\mathbf{r}}_{k-1}\}$ 可以相互线性表出, 故容易得到式 (4.249c) 和式 (4.249d). 证毕. \square

此外, 若定义

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{q}_k, \quad k = 0, 1, 2, \dots,$$

则迭代式 (4.248e) 中的向量 \mathbf{r}_k 刚好为 $\mathbf{r}_k = \mathbf{b} - \mathbf{A} \mathbf{x}_k$, 且由定理 4.16 的 (1) 和 (4) 可知

$$\mathbf{r}_k \perp \mathcal{K}_k(\mathbf{A}^T, \tilde{\mathbf{r}}_0),$$

这表明 \mathbf{x}_k 满足残量正交化条件式 (4.235). 综上所述, 就得到了如下算法.

算法 4.19 (BCG 方法) 给定线性方程组 (4.234), 初始向量 x_0 和 $\varepsilon > 0$. 本算法计算向量 x_k , 使得 $\|r_k\|_2/\|r_0\|_2 \leq \varepsilon$, 其中 $r_k = b - Ax_k$.

选取 x_0 ; 计算 $r_0 = b - Ax_0$; $q_{-1} = \tilde{q}_{-1} = 0$; $\rho_{-1} = 1$;

选择 \tilde{r}_0 满足 $\tilde{r}_0^T r_0 \neq 0$; $k = 0$;

while ($\|r_k\|_2/\|r_0\|_2 > \varepsilon$)

$\rho_k = \tilde{r}_k^T r_k$; $\beta_k = \rho_k/\rho_{k-1}$;

$q_k = r_k + \beta_k q_{k-1}$; $\tilde{q}_k = \tilde{r}_k + \beta_k \tilde{q}_{k-1}$;

$\sigma_k = \tilde{q}_k^T A q_k$; $\alpha_k = \rho_k/\sigma_k$;

$x_{k+1} = x_k + \alpha_k q_k$; $r_{k+1} = r_k - \alpha_k A q_k$;

$\tilde{r}_{k+1} = \tilde{r}_k - \alpha_k A^T \tilde{q}_k$; $k = k + 1$;

end

注 4.9 若计算过程中出现 $\rho_{k-1} = 0$ 或 $\sigma_k = 0$, 但还没有满足收敛性条件, 则此时算法就发生中断. 此外, 如果希望同时求解对偶方程 $A^T \tilde{x} = \tilde{b}$, 则只需在上述算法中增加两句 $\tilde{r}_0 = \tilde{b} - A^T \tilde{x}_0$ 和 $\tilde{x}_{k+1} = \tilde{x}_k + \alpha_k \tilde{q}_k$ 即可.

BCG 方法的 MATLAB 程序如下:

```
%BCG方法程序-bcg.m
function [x,k,time,res,resvec]=bcg(A,b,x,max_it,tol)
tic;n=length(b);
q=zeros(n,1);qt=q;rho=1;
r=b-A*x;mr=norm(r);k=0;
rt=ones(n,1);%rt的选取会影响收敛速度
%rt=r/mr; %rt=zeros(n,1);rt(1)=1;
while (k<max_it)
    k=k+1;
    res=norm(r)/mr;resvec(k)=res;
    if (res<tol), break; end
    rho1=rt'*r;beta=rho1/rho;
    q=r+beta*q;qt=rt+beta*qt;
    Aq=A*q; sigma=qt'*Aq;
    alpha=rho1/sigma;
    x=x+alpha*q;
    r=r-alpha*Aq;
    rt=rt-alpha*(A'*qt);
    rho=rho1;
end
time=toc;
```

例 4.15 假设线性方程组的系数矩阵 A 由 MATLAB 命令

$$\text{gallery('lotkin', n)}, \quad n = 1000,$$

产生, 这是一个 Hilbert 矩阵 (即其 (i, j) 元素为 $1/(i+j-1)$), 但其第 1 行的元素都换成了 1. 这个矩阵非常病态, 其条件数为 6.65×10^{21} . 再假设真解 x^* 是分量都为 1 的向量, 从而该方程组的右端项为 $b = Ax^*$. 将 BCG 算法应用到该线性方程组上, 取算法中的 $\tilde{r}_0 = \text{ones}(n, 1)$. 迭代在 60 步后收敛 ($\varepsilon = 10^{-10}$), 计算得到的近似解 \hat{x} 和真解 x^* 之间的绝对误差为

$$\|\hat{x} - x^*\|_2 = 0.6727,$$

但计算解 \hat{x} 的残量满足

$$\|b - A\hat{x}\|_2 = 2.6526 \times 10^{-8}.$$

迭代过程的收敛轨迹如图 4.16 所示, 其中横坐标为迭代步数 k , 纵坐标为 $\lg \|r_k\|_2$, 这里 r_k 是第 k 步得到的残量. 请注意观察, 与前面几个方法相比, 此时 $\|r_k\|_2/\|r_0\|_2$ 不再单调下降, 而是在振荡中下降着.

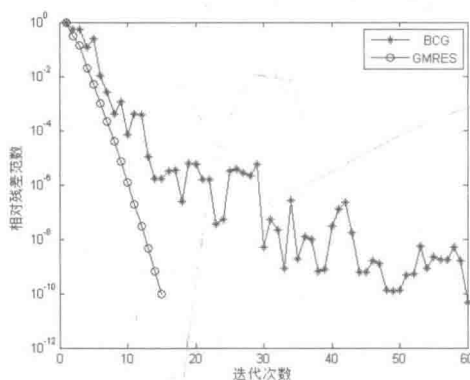


图 4.16 BCG 算法的收敛特性

4.8.2 CGS 方法

CGS 方法 (共轭梯度平方法, Conjugate Gradient Squared) 是为了避免 BCG 方法中 A^T 的出现而提出的.

仔细观察 BCG 的迭代格式就会发现, A^T 只在计算 \tilde{r}_k 时用到, 而这一量又在计算 ρ_k 和 σ_k 时用到. 再利用导出 BCG 迭代格式的多项式形式的迭代格式 (4.244) 立即可以看出, 这两个量的计算完全可以避开 A^T 的使用. 事实上, 有

$$\begin{cases} \rho_k = \langle \varphi_k, \varphi_k \rangle = \tilde{r}_0^T \varphi_k(A) \varphi_k(A) r_0 = \langle \varphi_0, \varphi_k^2 \rangle, \\ \sigma = \langle \psi_k, \vartheta \psi_k \rangle = \tilde{r}_0^T \psi_k(A) A \psi_k(A) r_0 = \langle \psi_0, \vartheta \psi_k^2 \rangle. \end{cases} \quad (4.252)$$

因此, 若定义

$$r_k = \varphi_k^2(A) r_0, \quad q_k = \psi_k^2(A) r_0, \quad (4.253)$$

则有

$$\rho_k = \tilde{\mathbf{r}}_0^T \mathbf{r}_k, \quad \sigma_k = \tilde{\mathbf{r}}_0^T \mathbf{A} \mathbf{q}_k. \quad (4.254)$$

这样一来, 要用式 (4.253) 来高效地计算 \mathbf{r}_k 和 \mathbf{q}_k , 关键是如何高效地计算 φ_k^2 和 ψ_k^2 . 由式 (4.244) 即可导出

$$\psi_k^2 = (\varphi_k + \beta_k \psi_{k-1})^2 = \varphi_k^2 + 2\beta_k \varphi_k \psi_{k-1} + \beta_k^2 \psi_{k-1}^2, \quad (4.255)$$

$$\varphi_{k+1}^2 = (\varphi_k - \alpha_k \vartheta \psi_k)^2 = \varphi_k^2 - 2\alpha_k \vartheta \varphi_k \psi_k + \alpha_k^2 \vartheta^2 \psi_k^2, \quad (4.256)$$

其中又涉及 $\varphi_k \psi_{k-1}$ 和 $\varphi_k \psi_k$. 再利用式 (4.244), 有

$$\varphi_{k+1} \psi_k = \varphi_k \psi_k - \alpha_k \vartheta \psi_k^2, \quad (4.257)$$

$$\varphi_k \psi_k = \varphi_k^2 + \beta_k \varphi_k \psi_{k-1}. \quad (4.258)$$

再定义

$$\mathbf{p}_k = \varphi_k(\mathbf{A}) \psi_{k-1}(\mathbf{A}) \mathbf{r}_0, \quad \mathbf{u}_k = \varphi_k(\mathbf{A}) \psi_k(\mathbf{A}) \mathbf{r}_0, \quad (4.259)$$

由式 (4.255)~式 (4.258) 得到如下递推公式:

$$\begin{cases} \mathbf{q}_k = \mathbf{r}_k + 2\beta_k \mathbf{p}_k + \beta_k^2 \mathbf{q}_{k-1}, \\ \mathbf{r}_{k+1} = \mathbf{r}_k - 2\alpha_k \mathbf{A} \mathbf{u}_k + \alpha_k^2 \mathbf{A}^2 \mathbf{q}_k, \\ \mathbf{p}_{k+1} = \mathbf{u}_k - \alpha_k \mathbf{A} \mathbf{q}_k, \\ \mathbf{u}_k = \mathbf{r}_k + \beta_k \mathbf{p}_k. \end{cases} \quad (4.260)$$

由式 (4.260) 的第 1 式和第 4 式, 得

$$\mathbf{q}_k = \mathbf{u}_k + \beta_k (\mathbf{p}_k + \beta_k \mathbf{q}_{k-1}), \quad (4.261)$$

再由式 (4.260) 的第 2 式和第 3 式, 得

$$\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{A} (\mathbf{u}_k + \mathbf{p}_{k+1}). \quad (4.262)$$

由式 (4.262) 可知, 若定义

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k (\mathbf{u}_k + \mathbf{p}_{k+1}), \quad (4.263)$$

则 $\mathbf{r}_k = \mathbf{b} - \mathbf{A} \mathbf{x}_k$ 正好是对应近似解的残量.

综合上面的讨论, 便有如下的共轭梯度平方法.

算法 4.20 (CGS 方法) 给定线性方程组 (4.234), 初始向量 \mathbf{x}_0 和容许误差 $\varepsilon > 0$.

本算法计算 \mathbf{x}_k , 使得 $\|\mathbf{r}_k\|_2 / \|\mathbf{r}_0\|_2 \leq \varepsilon$, 其中 $\mathbf{r}_k = \mathbf{b} - \mathbf{A} \mathbf{x}_k$.

选取 \mathbf{x}_0 ; 计算 $\mathbf{r}_0 = \mathbf{b} - \mathbf{A} \mathbf{x}_0$; $\mathbf{q}_{-1} = \mathbf{p}_0 = 0$; $\rho_{-1} = 1$;

选择 $\tilde{\mathbf{r}}_0$ 满足 $\tilde{\mathbf{r}}_0^T \mathbf{r}_0 \neq 0$; $k = 0$;

while ($\|\mathbf{r}_k\|_2 / \|\mathbf{r}_0\|_2 > \varepsilon$)

$\rho_k = \tilde{\mathbf{r}}_0^T \mathbf{r}_k$; $\beta_k = \rho_k / \rho_{k-1}$;

$$\begin{aligned}
\mathbf{u}_k &= \mathbf{r}_k + \beta_k \mathbf{p}_k; \\
\mathbf{q}_k &= \mathbf{u}_k + \beta_k (\mathbf{p}_k + \beta_k \mathbf{q}_{k-1}); \\
\mathbf{q}_k &= \mathbf{A} \mathbf{q}_k; \sigma_k = \tilde{\mathbf{r}}_0^T \mathbf{q}_k; \\
\alpha_k &= \rho_k / \sigma_k; \mathbf{p}_{k+1} = \mathbf{u}_k - \alpha_k \mathbf{q}_k; \\
\mathbf{z}_k &= \alpha_k (\mathbf{u}_k + \mathbf{p}_{k+1}); \\
\mathbf{x}_{k+1} &= \mathbf{x}_k + \mathbf{z}_k; \mathbf{r}_{k+1} = \mathbf{r}_k - \mathbf{A} \mathbf{z}_k; \\
k &= k + 1;
\end{aligned}$$

end

算法 4.20 确实避免了 \mathbf{A}^T 的出现, 每次迭代只需作两次系数矩阵 \mathbf{A} 与向量的乘积. 由前面的推导过程可知, 算法 4.20 产生的近似解的残量为 $\mathbf{r}_k^{\text{CGS}} = \varphi_k^2(\mathbf{A}) \mathbf{r}_0$, 而 BCG 算法产生的近似解的残量为 $\mathbf{r}_k^{\text{BCG}} = \varphi_k(\mathbf{A}) \mathbf{r}_0$. 当 $\mathbf{r}_k^{\text{BCG}}$ 趋向于零时, $\varphi_k(\mathbf{A})$ 就是一个收缩因子. 从这个意义上讲, $\mathbf{r}_k^{\text{CGS}}$ 的收敛到零的速度应该是 $\mathbf{r}_k^{\text{BCG}}$ 的两倍. 但这个算法的缺点是 $\|\mathbf{r}_k^{\text{CGS}}\|_2$ 随着 k 的增加会发生激烈的抖动.

例 4.16 仍然考虑例 4.15 中的 \mathbf{A} 和 \mathbf{b} . 将 CGS 算法应用到该线性方程组上, 取算法中的 $\tilde{\mathbf{r}}_0 = (1, 0, \dots, 0)^T$, 迭代在 142 步后收敛 ($\varepsilon = 10^{-10}$), 计算得到的近似解 $\hat{\mathbf{x}}$ 和真解 \mathbf{x}^* 之间的误差为

$$\|\hat{\mathbf{x}} - \mathbf{x}^*\|_2 = 0.8086,$$

但计算解 $\hat{\mathbf{x}}$ 的残量满足

$$\|\mathbf{b} - \mathbf{A}\hat{\mathbf{x}}\|_2 = 7.0413 \times 10^{-8}.$$

迭代过程的收敛轨迹如图 4.17 所示, 其中横坐标为迭代步数 k , 纵坐标为相对残差 $\|\mathbf{r}_k\|_2 / \|\mathbf{r}_0\|_2$, 这里 \mathbf{r}_k 是第 k 步得到的残差向量. 此时 $\|\mathbf{r}_k\|_2 / \|\mathbf{r}_0\|_2$ 与 BCG 方法一样也是在剧烈的振荡中艰难地下降着, 但注意到此例反映不出前面分析中 CGS 的收敛速度大约是 BCG 的两倍, 或许是受 $\tilde{\mathbf{r}}_0$ 选取的影响.

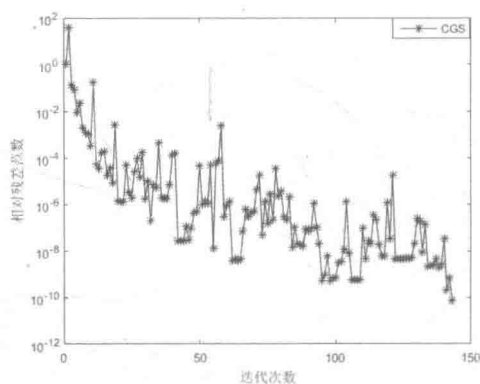


图 4.17 CGS 算法的收敛特性

4.8.3 BCGSTAB 方法

BCGSTAB 方法 (稳定化双共轭梯度法, Bi-Conjugate Gradient Stabilized) 是为了改进 CGS 方法之残量的范数剧烈抖动而提出的. 这一方法的基本思想是 CGS 方法的残量 $\mathbf{r}_k^{\text{CGS}}$ 满足

$$\mathbf{r}_k^{\text{CGS}} = \varphi_k(\mathbf{A})\mathbf{r}_k^{\text{BCG}} = \varphi_k^2(\mathbf{A})\mathbf{r}_0, \quad (4.264)$$

可以考虑不用多项式 φ_k , 而是选择一个其他的 k 次多项式 $\tilde{\varphi}_k$, 使

$$\mathbf{r}_k = \tilde{\varphi}_k(\mathbf{A})\mathbf{r}_k^{\text{BCG}} = \tilde{\varphi}_k(\mathbf{A})\varphi_k(\mathbf{A})\mathbf{r}_0, \quad (4.265)$$

以期待这样选取的相对残差范数 $\|\mathbf{r}_k\|_2/\|\mathbf{r}_0\|_2$ 的振荡性有所改进.

一种选择 $\tilde{\varphi}_k$ 的方法是根据递推公式

$$\tilde{\varphi}_0(t) = 1, \quad \tilde{\varphi}_{k+1}(t) = (1 - \omega_{k+1}t)\tilde{\varphi}_k(t), \quad (4.266)$$

式中: ω_{k+1} 为待定参数. BCGSTAB 方法正是利用这一参数的可选择性来改进其相对残差范数的振荡性的.

利用式 (4.244) 和式 (4.266), 可导出 $\tilde{\varphi}_{k+1}\varphi_{k+1}$ 的递推计算公式

$$\begin{aligned} \tilde{\varphi}_{k+1}\varphi_{k+1} &= (1 - \omega_{k+1}\vartheta)\tilde{\varphi}_k(\varphi_k - \alpha_k\vartheta\psi_k) \\ &= (1 - \omega_{k+1}\vartheta)(\tilde{\varphi}_k\varphi_k - \alpha_k\vartheta\tilde{\varphi}_k\psi_k), \end{aligned} \quad (4.267)$$

$$\begin{aligned} \tilde{\varphi}_k\psi_k &= \tilde{\varphi}_k(\varphi_k + \beta_k\psi_{k-1}) \\ &= \tilde{\varphi}_k\varphi_k + \beta_k(1 - \omega_k\vartheta)\tilde{\varphi}_{k-1}\psi_{k-1}. \end{aligned} \quad (4.268)$$

现定义

$$\mathbf{r}_k = \tilde{\varphi}_k(\mathbf{A})\varphi_k(\mathbf{A})\mathbf{r}_0, \quad \mathbf{p}_k = \tilde{\varphi}_k(\mathbf{A})\psi_k(\mathbf{A})\mathbf{r}_0, \quad (4.269)$$

则由式 (4.267) 和式 (4.268), 得

$$\mathbf{r}_{k+1} = (\mathbf{I} - \omega_{k+1}\mathbf{A})(\mathbf{r}_k - \alpha_k\mathbf{A}\mathbf{p}_k), \quad (4.270)$$

$$\mathbf{p}_k = \mathbf{r}_k + \beta_k(\mathbf{I} - \omega_k\mathbf{A})\mathbf{p}_{k-1}. \quad (4.271)$$

下面来考虑 α_k 和 β_k 的计算问题. 由式 (4.244) 可知

$$\alpha_k = \frac{\rho_k}{\sigma_k}, \quad \beta_k = \frac{\rho_k}{\rho_{k-1}}, \quad (4.272)$$

式中:

$$\rho_k = \langle \varphi_k, \varphi_k \rangle, \quad \sigma_k = \langle \psi_k, \vartheta\psi_k \rangle. \quad (4.273)$$

由定理 4.16 可知

$$\begin{aligned} \varphi_k(\mathbf{A})\mathbf{r}_0 &= \mathbf{r}_k^{\text{BCG}} \perp \mathcal{K}_k(\mathbf{A}^T, \tilde{\mathbf{r}}_0), \\ \psi_k(\mathbf{A})\mathbf{r}_0 &= \mathbf{q}_k^{\text{BCG}} \perp \mathbf{A}^T\mathcal{K}_k(\mathbf{A}^T, \tilde{\mathbf{r}}_0). \end{aligned}$$

于是对任意的 $\psi \in \mathcal{P}_{k-1}$, 有

$$\langle \varphi_k, \psi \rangle = \tilde{\mathbf{r}}_0^T \psi(\mathbf{A}) \varphi_k(\mathbf{A}) \mathbf{r}_0 = (\psi(\mathbf{A}^T) \tilde{\mathbf{r}}_0)^T (\varphi_k(\mathbf{A}) \mathbf{r}_0) = 0, \quad (4.274)$$

$$\langle \psi_k, \vartheta \psi \rangle = \tilde{\mathbf{r}}_0^T \psi(\mathbf{A}) \mathbf{A} \psi_k(\mathbf{A}) \mathbf{r}_0 = (\mathbf{A}^T \psi(\mathbf{A}^T) \tilde{\mathbf{r}}_0)^T (\psi_k(\mathbf{A}) \mathbf{r}_0) = 0. \quad (4.275)$$

设 φ_k 和 $\tilde{\varphi}_k$ 的首项系数分别为 ξ_k 和 η_k , 则由式 (4.244) 和式 (4.266) 可知, 它们可递推地计算

$$\xi_{k+1} = -\alpha_k \xi_k, \quad \eta_{k+1} = -\omega_{k+1} \eta_k, \quad (4.276)$$

其中 $\xi_0 = \eta_0 = 1$, 这里用到了

$$\psi_k = \varphi_k + \beta_k \psi_{k-1}$$

蕴涵着 ψ_k 和 φ_k 有相同的首项系数. 这样

$$\psi = \psi_k - \frac{\xi_k}{\eta_k} \tilde{\varphi}_k \quad \text{和} \quad \varphi = \varphi_k - \frac{\xi_k}{\eta_k} \tilde{\varphi}_k$$

均为 $k-1$ 次多项式, 从而利用式 (4.274) 和式 (4.275) 有

$$\begin{aligned} \rho_k &= \langle \varphi_k, \varphi_k \rangle = \langle \varphi_k, \varphi + \frac{\xi_k}{\eta_k} \tilde{\varphi}_k \rangle = \langle \varphi_k, \frac{\xi_k}{\eta_k} \tilde{\varphi}_k \rangle \\ &= \frac{\xi_k}{\eta_k} \tilde{\mathbf{r}}_0^T \tilde{\varphi}_k(\mathbf{A}) \varphi_k(\mathbf{A}) \mathbf{r}_0 = \frac{\xi_k}{\eta_k} \tilde{\mathbf{r}}_0^T \mathbf{r}_k, \end{aligned} \quad (4.277)$$

$$\begin{aligned} \sigma_k &= \langle \psi_k, \vartheta \psi_k \rangle = \langle \psi + \frac{\xi_k}{\eta_k} \tilde{\varphi}_k, \vartheta \psi_k \rangle = \langle \frac{\xi_k}{\eta_k} \tilde{\varphi}_k, \vartheta \psi_k \rangle \\ &= \frac{\xi_k}{\eta_k} \tilde{\mathbf{r}}_0^T \mathbf{A} \tilde{\varphi}_k(\mathbf{A}) \psi_k(\mathbf{A}) \mathbf{r}_0 = \frac{\xi_k}{\eta_k} \tilde{\mathbf{r}}_0^T \mathbf{A} \mathbf{p}_k. \end{aligned} \quad (4.278)$$

因此,

$$\alpha_k = \frac{\rho_k}{\sigma_k} = \frac{\tilde{\mathbf{r}}_0^T \mathbf{r}_k}{\tilde{\mathbf{r}}_0^T \mathbf{A} \mathbf{p}_k}, \quad (4.279)$$

$$\beta_k = \frac{\rho_k}{\rho_{k-1}} = \frac{\xi_k}{\xi_{k-1}} \frac{\eta_{k-1}}{\eta_k} \frac{\tilde{\mathbf{r}}_0^T \mathbf{r}_k}{\tilde{\mathbf{r}}_0^T \mathbf{r}_{k-1}} = \frac{\alpha_{k-1}}{\omega_k} \frac{\tilde{\mathbf{r}}_0^T \mathbf{r}_k}{\tilde{\mathbf{r}}_0^T \mathbf{r}_{k-1}}, \quad (4.280)$$

其中最后一个等式用到了式 (4.276).

到此为止 ω_{k+1} 仍然是自由的, 下面确定 ω_{k+1} . 为了符号简单起见, 记

$$\mathbf{s}_k = \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k, \quad (4.281)$$

有

$$\mathbf{r}_{k+1} = (\mathbf{I} - \omega_{k+1} \mathbf{A}) \mathbf{s}_k. \quad (4.282)$$

自然选择 ω_{k+1} 使

$$\|\mathbf{r}_{k+1}\|_2 = \min_{\omega} \|(\mathbf{I} - \omega \mathbf{A}) \mathbf{s}_k\|_2.$$

解此最小二乘问题, 得

$$\omega_{k+1} = \frac{\mathbf{s}_k^T \mathbf{A} \mathbf{s}_k}{\mathbf{s}_k^T \mathbf{A}^T \mathbf{A} \mathbf{s}_k}. \quad (4.283)$$

正是 ω_{k+1} 的这种极小化取法 $\|r_{k+1}\|_2$ 而使得的振荡性有所改善.

注意: 式 (4.282) 又可写为

$$r_{k+1} = s_k - \omega_{k+1} A s_k = r_k - \alpha_k A p_k - \omega_{k+1} A s_k, \quad (4.284)$$

故可定义

$$x_{k+1} = x_k + \alpha_k p_k + \omega_{k+1} s_k. \quad (4.285)$$

这样定义 x_{k+1} 后, 其残量正好是式 (4.282) 所定义的 r_{k+1} .

综合上面的讨论, 便得到如下的稳定化双共轭梯度算法.

算法 4.21 (BCGSTAB 方法) 给定线性方程组 (4.234), 初始向量 x_0 和 $\varepsilon > 0$. 本算法计算 x_k , 使得 $\|r_k\|_2 / \|r_0\|_2 \leq \varepsilon$, 其中 $r_k = b - A x_k$.

选取 x_0 ; 计算 $r_0 = b - A x_0$; $p_0 = r_0$;

选择 \tilde{r}_0 使 $\rho_0 = \tilde{r}_0^T r_0 \neq 0$; $k = 0$;

while ($\|r_k\|_2 > \|r_0\|_2 \varepsilon$)

$k = k + 1$;

$u_k = A p_k$; $\sigma_k = \tilde{r}_0^T u_k$;

$\alpha_k = \rho_k / \sigma_k$; $s_k = r_k - \alpha_k u_k$;

$q_k = A s_k$; $\omega_{k+1} = (s_k^T q_k) / (q_k^T q_k)$;

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

$r_{k+1} = s_k - \omega_{k+1} q_k$;

$\rho_{k+1} = \tilde{r}_0^T r_{k+1}$;

$\beta_{k+1} = (\alpha_k \rho_{k+1}) / (\omega_{k+1} \rho_k)$;

$p_{k+1} = r_{k+1} + \beta_{k+1} (p_k - \omega_{k+1} u_k)$;

end

算法 4.21 也避免了 A^T 的出现, 而且每次迭代也只需作两次矩阵向量乘法. 实际应用显示这种方法收敛得较快, 而且残量的范数较平稳, 是一种比较好的算法.

BCGSTAB 方法的 MATLAB 程序如下:

```
%BCGSTAB方法程序-bcgstab.m
function [x,k,time,res,resvec]=bcgstab(A,b,x,max_it,tol)
tic; r=b-A*x; p=r; mr=norm(r);
rt=ones(length(b),1); %rt=r;
%rt=zeros(length(b),1); rt(1)=1;
rho=rt'*r; k=0;
while (k<=max_it)
    k=k+1; u=A*p; sigma=rt'*u;
```

```

alpha=rho/sigma; s=r-alpha*u;
v=A*s; omega=(s'*v)/(v'*v);
x=x+alpha*p+omega*s;
r=s-omega*v; rho1=rt'*r;
beta=(alpha*rho1)/(omega*rho);
p=r+beta*(p-omega*u);
res=norm(r)/mr;resvec(k)=res;
if (res<tol), break; end
rho=rho1;
end
time=toc;

```

例 4.17 仍然考虑例 4.15 中的 A 和 b . 将 BCGSTAB 算法应用到该线性方程组上, 取算法中的 $\tilde{r}_0 = \text{ones}(n, 1)$, 迭代在 101 步后收敛 ($\epsilon = 10^{-10}$), 计算得到的近似解 \hat{x} 和真解 x^* 之间的绝对误差为

$$\|\hat{x} - x^*\|_2 = 0.9536,$$

但计算解 \hat{x} 的残量满足

$$\|b - A\hat{x}\|_2 = 7.8561 \times 10^{-8}.$$

迭代过程的收敛轨迹如图 4.18 所示, 其中横坐标为迭代步数 k , 纵坐标为相对残差 $\|r_k\|_2/\|r_0\|_2$, 这里 r_k 是第 k 步得到的残差向量. 此时的收敛速度与 CGS 方法基本一样, 但是注意 $\|r_k\|_2/\|r_0\|_2$ 已经变为呈平稳下降.

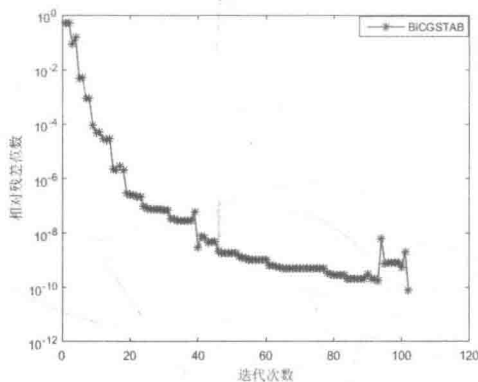


图 4.18 BCGSTAB 算法的收敛特性

最后, 关于预处理 BCGSTAB 算法, 对预处理方程组

$$\hat{A}\hat{x} = \hat{b}, \quad (4.286)$$

式中:

$$\hat{A} = M_1^{-1} A M_2^{-1}, \quad \hat{x} = M_2 x, \quad \hat{b} = M_1^{-1} b.$$

对预处理方程组 (4.286) 应用 BCGSTAB 算法, 再作代换:

$$\begin{aligned} \hat{u}_k &= M_1^{-1} u_k, \quad \hat{s}_k = M_1^{-1} s_k, \quad \hat{q}_k = M_1^{-1} q_k, \\ \hat{r}_k &= M_1^{-1} r_k, \quad \hat{p}_k = M_1^{-1} p_k, \quad \hat{x}_k = M_2 x_k, \quad \hat{\tilde{r}}_0 = M_1^T \tilde{r}_0. \end{aligned}$$

便得到如下预处理稳定化双共轭梯度算法.

算法 4.22 (PBCGSTAB 方法) 给定线性方程组 (4.234), 初始向量 x_0 和 $\varepsilon > 0$ 及预处理矩阵 $M = M_1 M_2$. 本算法计算 x_k , 使得 $\|r_k\|_2 / \|r_0\|_2 \leq \varepsilon$, 其中 $r_k = b - Ax_k$.

选取 x_0 ; 计算 $r_0 = b - Ax_0$; $p_0 = r_0$;

选择 \tilde{r}_0 使 $\rho_0 = \tilde{r}_0^T r_0 \neq 0$; $k = 0$;

while ($\|r_k\|_2 / \|r_0\|_2 > \varepsilon$)

$k = k + 1$;

由 $My_k = p_k$ 解得 y_k ;

$u_k = Ay_k$; $\sigma_k = \tilde{r}_0^T u_k$;

$\alpha_k = \rho_k / \sigma_k$; $s_k = r_k - \alpha_k u_k$;

由 $Mz_k = s_k$ 解得 z_k ; $q_k = Az_k$;

由 $M_1[\xi_k, \eta_k] = [q_k, s_k]$ 解得 $[\xi_k, \eta_k]$;

$\omega_{k+1} = (\xi_k^T \eta_k) / (\xi_k^T \xi_k)$;

$x_{k+1} = x_k + \alpha_k y_k + \omega_{k+1} z_k$;

$r_{k+1} = s_k - \omega_{k+1} q_k$;

$\rho_{k+1} = \tilde{r}_0^T r_{k+1}$; $\beta_{k+1} = (\alpha_k \rho_{k+1}) / (\omega_{k+1} \rho_k)$;

$p_{k+1} = r_{k+1} + \beta_{k+1} (p_k - \omega_{k+1} u_k)$;

end

若考虑例 4.15 中的 A 和 b . 将 PBCGSTAB 方法应用到该线性方程组上, 取算法中的 $\tilde{r}_0 = r_0$, 预处理矩阵 $M = M_1 M_2$ 为矩阵 A 的不完全 LU 分解, 则只需迭代 1 步即满足终止准则 ($\varepsilon = 10^{-10}$), 计算得到的近似解 \hat{x} 和真解 x^* 之间的误差为

$$\|\hat{x} - x^*\|_2 = 5.7303 \times 10^{-15},$$

计算解 \hat{x} 的残量满足

$$\|b - A\hat{x}\|_2 = 6.0738 \times 10^{-15}.$$

习题 4

4.1 设 $A \in \mathbb{R}^{n \times n}$ 是对称正定矩阵, p_1, \dots, p_k 是相互共轭正交的向量组, 即满足 $p_i^T A p_j = 0$ ($i \neq j$). 试证明: p_1, \dots, p_k 是线性无关的.

4.2 设 $A \in \mathbb{R}^{n \times n}$ 是对称正定矩阵, \mathcal{X} 是 \mathbb{R}^n 的一个 k 维子空间. 试证明: $x_k \in \mathcal{X}$ 满足

$$\|x_k - A^{-1}b\|_A = \min_{x \in \mathcal{X}} \|x - A^{-1}b\|_A$$

的充分必要条件是 $r_k = b - Ax_k$ 垂直于子空间 \mathcal{X} , 其中 $b \in \mathbb{R}^n$ 是任意给定的向量.

4.3 设对称正定矩阵 $A \in \mathbb{R}^{n \times n}$ 至多有 l 个互不相同的特征值. 试证明: 在没有舍入误差的前提下, 共轭梯度法至多 l 步即可得到方程组 $Ax = b$ 的精确解.

4.4 设对称矩阵 $A \in \mathbb{R}^{n \times n}$ 只有 l 个互不相同的特征值, $b \in \mathbb{R}^n$ 是任一向量, 试证明: 子空间 $\mathcal{X} = \text{span}\{b, Ab, \dots, A^{n-1}b\}$ 的维数至多为 l .

4.5 设 r_k 和 z_k 是由预处理共轭梯度法 (算法 4.2) 产生的. 证明: 若 $r_k \neq 0$, 则必有 $z_k^T r_k \neq 0$.

4.6 设

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix}, \quad b = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}.$$

取 $m = 2$, $x_0 = 0$. 分别用 GMRES 方法和 GMRES(m) 方法计算出 x_1 .

4.7 证明: 若 GMRES 方法中出现 $h_{k+1,k} = 0$, 则相应的近似解 x_k 为精确解. 反之亦然.

4.8 证明: BCG 方法的残量 r_k 也有三项递推公式.

4.9 设矩阵 $A \in \mathbb{R}^{n \times n}$ 非奇异. 证明 GMRES 方法中的最小二乘问题 (4.55) 是满秩的.

4.10 设矩阵 $A \in \mathbb{R}^{n \times n}$ 非奇异. 对于方程组 $Ax = b$, 使用内积 $(A^T A \cdot, \cdot)$, 导出在 Krylov 子空间

$$\text{span}\{r_0, Ar_0, A^2 r_0, \dots, A^k r_0\}$$

上的最佳近似解的计算公式.

第 5 章 线性最小二乘问题的数值解法

当给定的矩阵 $A \in \mathbb{C}^{n \times n}$ 非奇异时, 对任何向量 $b \in \mathbb{C}^n$, 线性方程组 $Ax = b$ 总有唯一解 $x = A^{-1}b$. 但在许多实际问题中, 矩阵 A 不是方阵, 甚至不是满秩的, 而且 $b \notin \mathcal{R}(A)$. 此时, 方程组 $Ax = b$ 可能无解或者有无穷多个解. 本章将讨论这样的一类线性方程组的有关理论及数值方法.

5.1 线性最小二乘问题的数学性质

本节讨论线性最小二乘问题的有关性质. 首先给出最小二乘问题的定义.

定义 5.1 设 $A \in \mathbb{C}^{m \times n}$, $b \in \mathbb{C}^m$, 确定 $x \in \mathbb{C}^n$ 使得

$$\|Ax - b\|_2 = \min_{z \in \mathbb{C}^n} \|Az - b\|_2. \quad (5.1)$$

问题 (5.1) 称为线性最小二乘问题 (Least Squares, LS 问题), 而 x 则称为最小二乘解或极小解. 称 $r(x) = b - Ax$ 为残差向量 (简称残量).

所有最小二乘解的集合记为 S_{LS} , 即

$$S_{LS} = \{x \in \mathbb{C}^n : x \text{ 满足 (5.1)}\}. \quad (5.2)$$

S_{LS} 中 2-范数最小者称为极小范数解, 记为 x_{LS} , 即

$$\|x_{LS}\|_2 = \min \{\|x\|_2 : x \in S_{LS}\}.$$

线性最小二乘问题 (5.1) 的解 x 又可称为线性方程组

$$Ax = b, \quad A \in \mathbb{C}^{m \times n}, \quad b \in \mathbb{C}^m \quad (5.3)$$

的最小二乘解, 即 x 在残量 $r(x) = b - Ax$ 的 2-范数最小的意义下满足方程组 (5.3). 当 $m > n$ 时称为超定方程组或矛盾方程组; 当 $m < n$ 时称为欠定方程组.

不难发现, 若将矩阵 A 写成 $A = [a_1, a_2, \dots, a_n]$, $a_i \in \mathbb{C}^m$, $i = 1, 2, \dots, n$, 则求解最小二乘问题 (5.1) 等价于求 $\{a_i\}_{i=1}^n$ 的线性组合使之与向量 b 之差的 2-范数达到最小. 可分为两种情况: 第一种是 $\{a_i\}_{i=1}^n$ 线性无关, 即 A 为列满秩; 第二种是 $\{a_i\}_{i=1}^n$ 线性相关, 即 A 为秩亏的. 下面分别针对这两种情形讨论最小二乘问题 (5.1) 极小解的数值解法.

5.1.1 最小二乘解的特征及一般表示

矩阵的广义逆是研究线性方程组 (5.3) 最小二乘解的一个重要而有力的工具. 下面讨论:

(1) 当方程组 (5.3) 有解时, 如何确定 $x_0 \in \mathbb{C}^n$, 使得

$$\|x_0\|_2 = \min_{Ax=b} \|x\|_2,$$

称这样的 \mathbf{x}_0 为方程组 (5.3) 的极小范数解.

(2) 当方程组 (5.3) 无解时, 如何确定 $\mathbf{x}_0 \in \mathbb{C}^n$, 使得

$$\|\mathbf{x}_0\|_2 = \min_{\min \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2} \|\mathbf{x}\|_2,$$

称这样的 \mathbf{x}_0 为方程组 (5.3) 的极小范数最小二乘解.

引理 5.1 设 $\mathbf{A} \in \mathbb{C}^{m \times n}$, 则有

(1) $[\mathcal{R}(\mathbf{A})]^\perp = \mathcal{N}(\mathbf{A}^H)$, 并且 $\mathbb{C}^m = \mathcal{R}(\mathbf{A}) \oplus \mathcal{N}(\mathbf{A}^H)$ (或 $\mathcal{R}(\mathbf{A}) = [\mathcal{N}(\mathbf{A}^H)]^\perp$).

(2) $[\mathcal{R}(\mathbf{A}^H)]^\perp = \mathcal{N}(\mathbf{A})$, 并且 $\mathbb{C}^n = \mathcal{R}(\mathbf{A}^H) \oplus \mathcal{N}(\mathbf{A})$ (或 $\mathcal{N}(\mathbf{A}) = [\mathcal{R}(\mathbf{A}^H)]^\perp$).

证明 (1) 将 \mathbf{A} 按列划分为 $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n]$, 其中 $\mathbf{a}_i \in \mathbb{C}^m$ ($i = 1, 2, \dots, n$), 由于

$$\mathcal{R}(\mathbf{A}) = \text{span}\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\},$$

所以

$$\begin{aligned} [\mathcal{R}(\mathbf{A})]^\perp &= \{\mathbf{x} : \mathbf{x} \perp \mathbf{a}_i, i = 1, 2, \dots, n\} \\ &= \{\mathbf{x} : \mathbf{a}_i^H \mathbf{x} = 0, i = 1, 2, \dots, n\} \\ &= \{\mathbf{x} : \mathbf{A}^H \mathbf{x} = \mathbf{0}\} = \mathcal{N}(\mathbf{A}^H), \end{aligned}$$

即 $\mathcal{R}(\mathbf{A}) \perp \mathcal{N}(\mathbf{A}^H)$, 从而 $\mathcal{R}(\mathbf{A}) + \mathcal{N}(\mathbf{A}^H)$ 为直和. 再由

$$\dim \mathcal{R}(\mathbf{A}) + \dim \mathcal{N}(\mathbf{A}^H) = \text{rank}(\mathbf{A}) + [m - \text{rank}(\mathbf{A})] = m,$$

可得 $\mathbb{C}^m = \mathcal{R}(\mathbf{A}) \oplus \mathcal{N}(\mathbf{A}^H)$.

(2) 在 (1) 中以 $\mathbf{A}^H \in \mathbb{C}^{n \times m}$ 代替 $\mathbf{A} \in \mathbb{C}^{m \times n}$ 即可得 (2). 证毕. □

设矩阵 \mathbf{A} 的奇异值分解为

$$\mathbf{A} = \mathbf{U} \begin{bmatrix} \boldsymbol{\Sigma}_r & \mathbf{O} \\ \mathbf{O} & \mathbf{O} \end{bmatrix} \mathbf{V}^H,$$

式中: $\mathbf{U} = [\mathbf{U}_1, \mathbf{U}_2]$, $\mathbf{V} = [\mathbf{V}_1, \mathbf{V}_2]$ 为酉矩阵, $\mathbf{U}_1 = [\mathbf{u}_1, \dots, \mathbf{u}_r]$, $\mathbf{U}_2 = [\mathbf{u}_{r+1}, \dots, \mathbf{u}_n]$, $\mathbf{V}_1 = [\mathbf{v}_1, \dots, \mathbf{v}_r]$, $\mathbf{V}_2 = [\mathbf{v}_{r+1}, \dots, \mathbf{v}_n]$; $\boldsymbol{\Sigma}_r = \text{diag}(\sigma_1, \dots, \sigma_r)$, $\sigma_1 \geq \dots \geq \sigma_r > 0$. 于是

$$\mathbf{A} = \mathbf{U}_1 \boldsymbol{\Sigma}_r \mathbf{V}_1^H, \quad \mathbf{A}^\dagger = \mathbf{V}_1 \boldsymbol{\Sigma}_r^{-1} \mathbf{U}_1^H, \quad \mathbf{A}^H = \mathbf{V}_1 \boldsymbol{\Sigma}_r \mathbf{U}_1^H. \quad (5.4)$$

容易验证

$$\mathbf{A} \mathbf{A}^\dagger = \mathbf{U}_1 \mathbf{U}_1^H, \quad \mathbf{A}^\dagger \mathbf{A} = \mathbf{V}_1 \mathbf{V}_1^H, \quad \mathbf{A}^H \mathbf{A} = \mathbf{V}_1 \boldsymbol{\Sigma}_r^2 \mathbf{V}_1^H. \quad (5.5)$$

由于 $\{\mathbf{u}_i\}_{i=1}^r$ 和 $\{\mathbf{v}_i\}_{i=r+1}^n$ 分别是 $\mathcal{R}(\mathbf{A})$ 和 $\mathcal{N}(\mathbf{A})$ 的标准正交基, 故 $\mathcal{R}(\mathbf{A})$ 和 $\mathcal{N}(\mathbf{A})$ 上的正交投影矩阵分别为:

$$\mathbf{P}_{\mathcal{R}(\mathbf{A})} = \mathbf{U}_1 \mathbf{U}_1^H = \mathbf{A} \mathbf{A}^\dagger, \quad \mathbf{P}_{\mathcal{N}(\mathbf{A})} = \mathbf{V}_2 \mathbf{V}_2^H = \mathbf{I}_n - \mathbf{V}_1 \mathbf{V}_1^H = \mathbf{I}_n - \mathbf{A}^\dagger \mathbf{A}. \quad (5.6)$$

而 $\mathcal{R}(A)^\perp = \mathcal{N}(A^H)$ 和 $\mathcal{N}(A)^\perp = \mathcal{R}(A^H)$ 上的正交投影矩阵分别为

$$P_{\mathcal{R}(A)^\perp} = I_m - AA^\dagger, \quad P_{\mathcal{N}(A)^\perp} = A^\dagger A. \quad (5.7)$$

有下面的定理.

定理 5.1 设 $A \in \mathbb{C}^{m \times n}$, $b \in \mathbb{C}^m$, 则线性方程组 (5.3) 有解的充要条件是

$$AA^\dagger b = b, \quad (5.8)$$

并且在有解时, 其通解为

$$x = A^\dagger b + (I - A^\dagger A)z, \quad (5.9)$$

其中 $z \in \mathbb{C}^n$ 任意.

证明 若 $AA^\dagger b = b$, 则显然方程组 (5.3) 有解 $x = A^\dagger b$. 反之, 若 $Ax = b$, 则

$$AA^\dagger b = AA^\dagger Ax = Ax = b.$$

下面证明其通解为式 (5.9). 事实上, 根据式 (5.6), 有

$$\mathcal{N}(A) = \{(I_n - A^\dagger A)z : z \in \mathbb{C}^n\}.$$

故线性方程组 (5.3) 的通解可以表示为 $x = A^\dagger b + (I - A^\dagger A)z$, 其中 $z \in \mathbb{C}^n$ 任意. 证毕. \square

注 5.1 式 (5.9) 表明: $x_0 = A^\dagger b$ 是方程组 $Ax = b$ 的一个解, 而 $(I - A^\dagger A)z$ 是对应齐次方程组 $Ax = 0$ 的通解.

定理 5.2 如果方程组 (5.3) 有解, 则它的极小范数解 x_0 唯一, 并且 $x_0 = A^\dagger b$.

证明 先证 $x_0 \in \mathcal{R}(A^H)$. 若 $x_0 \notin \mathcal{R}(A^H)$, 则由引理 5.1 (2) 作向量分解, 得

$$x_0 = y_0 + y_1, \quad y_0 \in \mathcal{R}(A^H), \quad y_1 \in \mathcal{N}(A) = [\mathcal{R}(A^H)]^\perp, \quad y_1 \neq 0.$$

由于 $y_0 \perp y_1$, 所以

$$\|x_0\|_2^2 = \|y_0\|_2^2 + \|y_1\|_2^2 > \|y_0\|_2^2.$$

但是 $Ay_0 = Ay_0 + Ay_1 = Ax_0 = b$, 所以 x_0 不是方程组 (5.3) 的极小范数解. 这与前提冲突, 故 $x_0 \in \mathcal{R}(A^H)$.

再证 x_0 唯一. 若 z_0 也是方程 (5.3) 的极小范数解, 则 $z_0 \in \mathcal{R}(A^H)$, 从而

$$x_0 - z_0 \in \mathcal{R}(A^H) = [\mathcal{N}(A)]^\perp.$$

另外, 由于 $A(x_0 - z_0) = 0$, 所以 $x_0 - z_0 \in \mathcal{N}(A)$, 故只能有 $x_0 - z_0 = 0$. 即 $x_0 = z_0$.

最后证 $x_0 = A^\dagger b$. 根据定理 5.1 可得方程组 (5.3) 的通解为

$$x = A^\dagger b + (I - A^\dagger A)z \quad (z \in \mathbb{C}^n \text{ 任意}).$$

取 $z = 0$, 则 $x_0 = A^\dagger b$ 是方程组 (5.3) 的一个解. 因为

$$\begin{aligned}(x_0, (I - A^\dagger A)z) &= z^H (I - A^\dagger A)^H x_0 \\ &= z^H (I - A^\dagger A) A^\dagger b \\ &= z^H (A^\dagger - A^\dagger A A^\dagger) b = 0,\end{aligned}$$

所以

$$\|x\|_2^2 = \|x_0\|_2^2 + \|(I - A^\dagger A)z\|_2^2 \geq \|x_0\|_2^2.$$

故 $x_0 = A^\dagger b$ 是方程组 (5.3) 的极小范数解. 证毕. \square

定理 5.3 如果线性方程组 (5.3) 无解, 则它的极小范数最小二乘解 x_0 唯一, 并且 $x_0 = A^\dagger b$.

证明 由于方程组 (5.3) 无解, 所以 $b \notin \mathcal{R}(A)$. 根据引理 5.1 (1) 作向量分解 $b = b_1 + b_2$, 其中 $b_1 \in \mathcal{R}(A)$, $b_2 \in \mathcal{N}(A^H) = [\mathcal{R}(A)]^\perp$ 且 $b_2 \neq 0$. 注意到 $Ax - b_1 \in \mathcal{R}(A)$ 及 $b_2 \in [\mathcal{R}(A)]^\perp$, 可得 $(Ax - b_1) \perp b_2$, 于是有

$$\|Ax - b\|_2^2 = \|(Ax - b_1) + (-b_2)\|_2^2 = \|Ax - b_1\|_2^2 + \|b_2\|_2^2,$$

这表明 $\min \|Ax - b\|_2$ 等价于 $\min \|Ax - b_1\|_2$.

因为 $b_1 \in \mathcal{R}(A)$, 所以方程组 $Ax = b_1$ 有解, 根据定理 5.2 可得它的唯一极小范数解为 $A^\dagger b_1$. 由于 $b_2 \in \mathcal{N}(A^H)$, 所以 $A^H b_2 = 0$. 于是得

$$A^\dagger b_2 = A^\dagger A A^\dagger b_2 = A^\dagger (A A^\dagger)^H b_2 = A^\dagger (A^\dagger)^H A^H b_2 = 0,$$

所以 $x_0 = A^\dagger b = A^\dagger (b_1 + b_2) = A^\dagger b_1$ 是方程组 $Ax = b_1$ 的唯一极小范数解, 也是方程组 (5.3) 的唯一极小范数最小二乘解. 证毕. \square

5.1.2 线性 LS 的等价性问题

1. 法方程

下面的定理给出了最小二乘问题极小解的一个刻画.

定理 5.4 x 是最小二乘问题 (5.1) 的极小解, 即 $x \in S_{LS}$ 的充分必要条件是 x 为方程

$$A^H A x = A^H b \quad (5.10)$$

的解, 其中式 (5.10) 称为最小二乘问题的法方程.

证明 注意到最小二乘问题 (5.1) 等价于极小化问题

$$\min \varphi(x) = \frac{1}{2} \|Ax - b\|_2^2 = \frac{1}{2} x^H (A^H A) x - (b^H A) x + \frac{1}{2} \|b\|_2^2.$$

由于 $A^H A \in \mathbb{C}^{n \times n}$ 是半正定矩阵, 因此 n 元实函数 $\varphi(x)$ 是凸函数, 故 x 是最小二乘问题 (5.1) 的极小解等价于

$$\nabla \varphi(x) = A^H(Ax - b) = 0.$$

证毕. □

定理 5.4 说明最小二乘问题 (5.1) 与法方程 (5.10) 是等价的, 即法方程 (5.10) 与最小二乘问题 (5.1) 有相同的解集.

2. KKT 方程

下面的定理给出最小二乘问题 (5.1) 的另一个等价性问题.

定理 5.5 设 $A \in \mathbb{C}^{m \times n}$, $b \in \mathbb{R}^n$. 则 x 和 $r = b - Ax$ 分别为最小二乘问题 (5.1) 的极小解和残量的充分必要条件是 x 和 r 为鞍点系统

$$\begin{bmatrix} I & A \\ A^H & O \end{bmatrix} \begin{bmatrix} r \\ x \end{bmatrix} = \begin{bmatrix} b \\ 0 \end{bmatrix} \quad (5.11)$$

的解. 上述线性系统称为最小二乘问题的 Karush-Kuhn-Tucker 方程 (KKT 方程).

证明 若 x 为最小二乘问题 (5.1) 的极小解, 而 $r = b - Ax$ 为其残量, 则

$$x = A^\dagger b + (I - A^\dagger A)z, \quad r = b - Ax = b - AA^\dagger b = (I - AA^\dagger)b.$$

由广义逆 A^\dagger 的性质及等式 $r + Ax = b$, 得

$$A^H r = A^H(I - AA^\dagger)b = [(I - AA^\dagger)A]^H b = 0.$$

故式 (5.11) 是相容的线性系统, 且 x 和 r 满足式 (5.11).

反之, 通过验证广义逆的四个条件, 可以验证

$$B^\dagger \equiv \begin{bmatrix} I & A \\ A^H & O \end{bmatrix}^\dagger = \begin{bmatrix} I - AA^\dagger & (A^\dagger)^H \\ A^\dagger & -A^\dagger(A^\dagger)^H \end{bmatrix}.$$

故式 (5.11) 的任一解向量 x, r 有如下形式

$$\begin{bmatrix} r \\ x \end{bmatrix} = B^\dagger \begin{bmatrix} b \\ 0 \end{bmatrix} + (I - B^\dagger B) \begin{bmatrix} y \\ z \end{bmatrix} = \begin{bmatrix} (I - AA^\dagger)b \\ A^\dagger b - (I - A^\dagger A)z \end{bmatrix},$$

式中: $y \in \mathbb{C}^m$, $z \in \mathbb{C}^n$ 为任意向量. 故满足式 (5.11) 的任一组向量 x, r 分别为式 (5.1) 的极小解和残量. 证毕. □

5.1.3 线性最小二乘问题的正则化

设矩阵 $A \in \mathbb{C}_r^{m \times n}$ 的奇异值分解为

$$A = U \Sigma V^H, \quad (5.12)$$

式中: $U = [u_1, u_2, \dots, u_m]$ 为 m 阶酉矩阵; $V = [v_1, v_2, \dots, v_n]$ 为 n 阶酉矩阵; $l = \min\{m, n\}$; $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0 = \sigma_{r+1} = \dots = \sigma_l = 0$,

$$\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_l). \quad (5.13)$$

则最小二乘问题 (5.1) 的极小范数最小二乘解 x_{LS} 可表示为

$$x_{LS} = A^\dagger b = \sum_{i=1}^r \frac{u_i^H b}{\sigma_i} v_i. \quad (5.14)$$

由于舍入误差的存在, 所有用于计算 A^\dagger 和 x_{LS} 的算法, 实际上是计算某个扰动矩阵 $\hat{A} = A + \Delta A$ 的广义逆 \hat{A}^\dagger 以及 $\hat{x}_{LS} = \hat{A}^\dagger \hat{b}$.

当 A 为病态的满秩矩阵, 即存在 $k (1 \leq k < l)$, 使得 $\sigma_k \gg \sigma_{k+1} \approx 0$. 由于在极小解 x_{LS} 的表达式中存在项

$$\frac{u_i^H b}{\sigma_i} v_i, \quad i = k+1, \dots, l,$$

因此, 对 A 的较小扰动, 将会使极小解 x_{LS} 产生很大的误差. 此时, 称 A 为数值秩亏的. 当矩阵 A 秩亏而 ΔA 很小, 则广义逆的不连续性表明了原来意义上的秩不再适用于数值计算. 当 \hat{A} 和 A 的秩不相同, \hat{A}^\dagger 和 A^\dagger , \hat{x}_{LS} 和 x_{LS} 可能会相差很大, 而且扰动 ΔA 越小, 其相差的程度会越大.

1. 截断的 LS 问题

最小二乘问题的第一种正则化方法是截断的 LS 问题. 首先给出矩阵 A 的 δ 秩的定义.

定义 5.2 设 $A \in \mathbb{C}^{m \times n}$ 和 $\delta > 0$ 给定. 称数

$$k = \min_{B \in \mathbb{C}^{m \times n}} \{\text{rank}(B) : \|A - B\|_2 \leq \delta\} \quad (5.15)$$

为矩阵 A 的 δ 秩.

由定义 5.2 和矩阵的降秩最佳逼近定理, 当 $k < l$ 时, 有

$$\|A - A_k\|_2 = \min_{\text{rank}(B) \leq k} \|A - B\|_2 = \sigma_{k+1},$$

式中:

$$A_k = \sum_{i=1}^k \sigma_i u_i v_i^H.$$

因此, 矩阵 A 的 δ 秩为 k 的充分必要条件是

$$\sigma_1 \geq \dots \geq \sigma_k > \sigma_{k+1} \geq \dots \geq \sigma_l.$$

对于最小二乘问题 (5.1), 截断的 LS 问题为

$$\|A_k x - b\|_2 = \min_{z \in \mathbb{C}^n} \|A_k z - b\|_2, \quad (5.16)$$

式中: \mathbf{A}_k 为 \mathbf{A} 的最佳秩- k 逼近. 此时, 式 (5.16) 的极小范数最小二乘解 $\bar{\mathbf{x}}_{\text{LS}}$ 可表示为

$$\bar{\mathbf{x}}_{\text{LS}} = \mathbf{A}_k^\dagger \mathbf{b} = \sum_{i=1}^k \frac{\mathbf{u}_i^H \mathbf{b}}{\sigma_i} \mathbf{v}_i. \quad (5.17)$$

2. Tikhonov 正则化

最小二乘问题的另一种正则化方法是 Tikhonov 正则化, 即考虑如下的正则化问题:

$$\|\mathbf{Ax} - \mathbf{b}\|_2^2 + \tau^2 \|\mathbf{Dx}\|_2^2 = \min_{\mathbf{z} \in \mathbb{C}^n} \|\mathbf{Az} - \mathbf{b}\|_2^2 + \tau^2 \|\mathbf{Dz}\|_2^2, \quad (5.18)$$

式中: $\tau > 0$; $\mathbf{D} = \text{diag}(d_1, d_2, \dots, d_n)$ 为正定的对角矩阵.

易见, 式 (5.18) 等价于

$$\left\| \begin{bmatrix} \tau \mathbf{D} \\ \mathbf{A} \end{bmatrix} \mathbf{x} - \begin{bmatrix} \mathbf{0} \\ \mathbf{b} \end{bmatrix} \right\|_2^2 = \min_{\mathbf{z} \in \mathbb{C}^n} \left\| \begin{bmatrix} \tau \mathbf{D} \\ \mathbf{A} \end{bmatrix} \mathbf{z} - \begin{bmatrix} \mathbf{0} \\ \mathbf{b} \end{bmatrix} \right\|_2^2, \quad (5.19)$$

当 $\tau > 0$ 时, 系数矩阵是列满秩的, 因此有唯一的最小二乘解.

设 $m \geq n$. 则当 $\mathbf{D} = \mathbf{I}$ 时, 式 (5.19) 系数矩阵的奇异值为 $\tilde{\sigma}_i = \sqrt{\sigma_i^2 + \tau^2}$, $i = 1, 2, \dots, n$. 此时式 (5.18) 的解可表示为

$$\mathbf{x}(\tau) = \sum_{i=1}^n \frac{(\mathbf{u}_i^H \mathbf{b}) \sigma_i}{\sigma_i^2 + \tau^2} \mathbf{v}_i = \sum_{i=1}^n \frac{(\mathbf{u}_i^H \mathbf{b}) \eta_i}{\sigma_i} \mathbf{v}_i, \quad \eta_i = \frac{\sigma_i^2}{\sigma_i^2 + \tau^2}, \quad (5.20)$$

式中: η_i 为滤波因子. 当 $\tau \ll \sigma_i$ 时, $\eta_i \approx 1$; 当 $\tau \gg \sigma_i$ 时, $\eta_i \ll 1$.

正则化问题 (5.18) 的优点是它的解可以通过 QR 分解

$$\begin{bmatrix} \tau \mathbf{D} \\ \mathbf{A} \end{bmatrix} = \mathbf{Q} \begin{bmatrix} \mathbf{R} \\ \mathbf{O} \end{bmatrix}$$

得到.

注 5.2 最小二乘问题的正则化已经把原来的最小二乘问题化为新的最小二乘问题. 事实上, 由式 (5.17) 和式 (5.20) 可以看出, 截断的 LS 问题和 Tikhonov 正则化的 LS 问题, 与最小二乘问题 (5.1) 的极小范数最小二乘解以及相应的残量各不相同. 正则化方法通常用于处理病态最小二乘问题和反问题等不适定问题.

5.2 求解满秩最小二乘问题的数值方法

本节假定最小二乘问题 (5.1) 中的矩阵 $\mathbf{A} \in \mathbb{C}^{m \times n}$ ($m \geq n$) 为列满秩矩阵, 即 $\text{rank}(\mathbf{A}) = n$. 此时最小二乘问题 (5.1) 有唯一的极小范数最小二乘解 \mathbf{x}_{LS} , 并且连续地依赖给定的数据 \mathbf{A} 和 \mathbf{b} . 对于这类问题, 介绍两种最基本的数值方法.

5.2.1 法方程方法

由定理 5.4, 最小二乘问题 (5.1) 等价于其法方程

$$A^H A x = A^H b.$$

当 $\text{rank}(A) = n$ 时, $A^H A$ 为 Hermite 正定矩阵, 法方程 (5.10) 的唯一解可以用 Cholesky 分解法求得. 写出算法步骤如下:

算法 5.1 (法方程 Cholesky 分解法) 给定 $A \in \mathbb{C}^{m \times n}$ ($m \geq n$) 为列满秩矩阵, $b \in \mathbb{C}^m$. 本算法计算 $\|Ax - b\|_2$ 的极小解 x_{LS} .

步 1, 对 n 阶 Hermite 正定矩阵 $B = A^H A$ 作 Cholesky 分解 $B = LL^H$, 其中 L 为下三角矩阵.

步 2, 依次解 $Ly = A^H b$, $L^H x = y$ 得到最小二乘问题 (5.1) 的解 x_{LS} .

n 阶 Hermite 正定矩阵 B 的 Cholesky 分解需要 $O(n^3/3)$ 次乘法. 而计算 $A^H A$ 需要 mn^2 次乘法. 在算法 5.1 中第 2 步的计算量 (约 $O(n^2)$ 次乘法) 与第 1 步相比可以忽略. 所以该算法需要的乘法次数为 $O(mn^2 + n^3/3)$. 当 $m \gg n$ 时, 算法的主要工作量在于 $A^H A$ 的计算.

5.2.2 QR 分解方法

本节考虑最小二乘问题 (5.1) 中的矩阵 $A \in \mathbb{R}^{m \times n}$ ($m \geq n$) 的情形. 根据正交矩阵保持向量 2-范数不变的性质, 对于任意的正交矩阵 $Q \in \mathbb{R}^{m \times m}$, 最小二乘问题 (5.1) 等价于

$$\|Q^T(Ax - b)\|_2 = \min \left\{ \|Q^T(Az - b)\|_2 : z \in \mathbb{R}^n \right\}. \quad (5.21)$$

这样, 就可望通过适当选取正交矩阵 Q , 使原问题 (5.1) 转化为较为容易求解的最小二乘问题 (5.21), 这就是正交化方法—QR 分解方法的基本思想.

设 A 有 QR 分解:

$$A = Q \begin{bmatrix} R \\ O \end{bmatrix} = Q_1 R, \quad (5.22)$$

式中: $Q \in \mathbb{R}^{m \times m}$ 为正交矩阵; Q_1 为 Q 的前 n 列组成的矩阵; $R \in \mathbb{R}^{n \times n}$ 为对角元均为正的上三角矩阵.

现取式 (5.21) 中的正交矩阵为分解式 (5.22) 中的 Q , 并记

$$f = Q^T b = \begin{bmatrix} Q_1^T \\ Q_2^T \end{bmatrix} b = \begin{bmatrix} f_1 \\ f_2 \end{bmatrix} \begin{matrix} n \\ m-n \end{matrix}, \quad (5.23)$$

则对任意的 $x \in \mathbb{C}^n$, 有

$$\begin{aligned} \|Q^T(Ax - b)\|_2^2 &= \left\| \begin{bmatrix} R \\ O \end{bmatrix} x - Q^T b \right\|_2^2 = \left\| \begin{bmatrix} R \\ O \end{bmatrix} x - \begin{bmatrix} f_1 \\ f_2 \end{bmatrix} \right\|_2^2 \\ &= \|Rx - f_1\|_2^2 + \|f_2\|_2^2. \end{aligned}$$

由此可知, \mathbf{x} 是最小二乘问题 (5.1) 的解, 当且仅当 \mathbf{x} 是上三角形方程组 $\mathbf{R}\mathbf{x} = \mathbf{f}_1$ 的解.

根据上述讨论, 有如下的 QR 分解算法. 在实施过程中, 采取对增广矩阵 $\tilde{\mathbf{A}} = [\mathbf{A}, \mathbf{b}]$ 进行 QR 分解:

$$\mathbf{H}_n \cdots \mathbf{H}_2 \mathbf{H}_1 [\mathbf{A}, \mathbf{b}] = [\mathbf{R}, \tilde{\mathbf{b}}] \implies \mathbf{Q}^T [\mathbf{A}, \mathbf{b}] = [\mathbf{R}, \tilde{\mathbf{b}}], \quad (5.24)$$

得

$$\mathbf{A} = \mathbf{Q}\mathbf{R}, \quad \tilde{\mathbf{b}} = \mathbf{Q}^T \mathbf{b},$$

式中: $\mathbf{Q} = \mathbf{H}_1 \mathbf{H}_2 \cdots \mathbf{H}_n$. 式 (5.24) 表明, 对增广矩阵 $\tilde{\mathbf{A}} = [\mathbf{A}, \mathbf{b}]$ 实施 QR 分解, 当把矩阵 \mathbf{A} 约化为上三角矩阵 \mathbf{R} 时, 向量 \mathbf{b} 约化成了 $\mathbf{Q}^T \mathbf{b}$, 这正是所需要的.

算法 5.2 (LS 问题 QR 分解法) 给定 $\mathbf{A} \in \mathbb{R}^{m \times n}$ ($m \geq n$) 为列满秩矩阵, $\mathbf{b} \in \mathbb{R}^m$. 本算法计算 $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$ 的极小解 \mathbf{x}_{LS} .

步 1, 计算增广矩阵 $\tilde{\mathbf{A}} = [\mathbf{A}, \mathbf{b}]$ 的 QR 分解并覆盖 $\tilde{\mathbf{A}}$.

步 2, 用回代法求解上三角形方程组 $\text{triu}(\tilde{\mathbf{A}}(1:n, 1:n))\mathbf{x} = \tilde{\mathbf{A}}(1:n, n+1)$, 得到最小二乘解 \mathbf{x}_{LS} , 其中记号 $\text{triu}(\mathbf{B})$ 表示提取矩阵 \mathbf{B} 的上三角部分组成的上三角形矩阵.

用算法 5.2 求解满秩最小二乘问题需要 $2mn^2 - 2n^3/3$ 个 flop. 更新 \mathbf{b} 所需的 $O(mn)$ 个 flop 和回代法所需的 $O(n^2)$ 个 flop 与分解 \mathbf{A} 所需的工作量相比是微不足道的.

注 5.3 算法 5.2 的主要工作量集中在增广矩阵 $\tilde{\mathbf{A}}$ 的 QR 分解. 可用如下三种方法之一实现这一分解:

- (1) Householder 方法.
- (2) Givens 正交化方法.
- (3) 修正的 Gram-Schmit 方法.

根据算法 5.2, 采用 Householder 变换 QR 分解编制 MATLAB 程序如下:

```
function [x]=ls_houseqr(A,b)
%用Householder变换QR分解求最小二乘问题min||Ax-b||
[m,n]=size(A);
[A]=house_qr([A,b]); %调用Householder变换QR分解程序
x=triu(A(1:n,1:n))\A(1:n,n+1);
```

例 5.1 用 MATLAB 程序 ls_houseqr.m 求解最小二乘问题 $\min \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$, 其中

$$\mathbf{A} = \begin{bmatrix} 2 & 3 & 4 & 5 \\ 4 & 3 & 2 & 1 \\ 4 & 5 & 6 & 7 \\ 9 & 5 & 7 & 2 \\ 4 & 2 & 5 & 3 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 20 \\ 22 \\ 35 \\ 42 \\ 50 \end{bmatrix}.$$

解 在 MATLAB 命令窗口输入:

```
>> A=[2 3 4 5; 4 3 2 1; 4 5 6 7; 9 5 7 2; 4 2 5 3];  
>> b=[20 22 35 42 50]';  
>> x=ls_houseqr(A,b)
```

即可得计算结果.

5.3 求解秩亏最小二乘问题的数值解法

本节讨论秩亏最小二乘问题 $\min \|Ax - b\|_2$ 的数值解法. 如果 $A \in \mathbb{R}_r^{m \times n}$ ($m \geq n$) 是秩亏损的, 即 $\text{rank}(A) = r < n$, 此时, 最小二乘问题 (5.1) 有无穷多个解, 且 5.2 节所介绍的处理满秩最小二乘问题的法方程方法和 QR 分解法都不再有效. 因此, 本节专门讨论秩亏最小二乘问题的数值求解方法.

5.3.1 列主元 QR 分解法

当矩阵 A 秩亏损时, 其 QR 分解不一定能给出列空间 $\mathcal{R}(A)$ 的一组标准正交基. 此时可计算经过列置换之后的矩阵 $\tilde{A} = AP$ 的 QR 分解来解决, 即 $AP = QR$, 其中 P 是置换矩阵.

下面介绍列主元 QR 分解法来解决秩亏的最小二乘问题. 设最小二乘问题 (5.1) 中的已知向量 b 分解为

$$b = b_1 + b_2, \quad b_1 \in \mathcal{R}(A), \quad b_2 \in \mathcal{R}(A)^\perp. \quad (5.25)$$

易证问题 (5.1) 等价于

$$Ax = b_1. \quad (5.26)$$

现假定 $\mathcal{R}(A) = \mathcal{R}(Q_1)$, 其中 $Q_1 \in \mathbb{R}^{m \times r}$ 且 $Q_1^T Q_1 = I_r$, 即 Q_1 的列构成 $\mathcal{R}(A)$ 的一组标准正交基. 则存在矩阵 $S \in \mathbb{R}^{r \times n}$ 和向量 $h \in \mathbb{R}^r$, 使得

$$A = Q_1 S, \quad b_1 = Q_1 h. \quad (5.27)$$

将式 (5.27) 代入式 (5.26), 并注意到 Q_1 的列线性无关, 可知 (5.26) 等价于

$$Sx = h. \quad (5.28)$$

显然方程组 (5.28) 总是有解的. 进一步, 由式 (5.27) 可知

$$S = Q_1^T A, \quad h = Q_1^T b_1 = Q_1^T (b - b_2) = Q_1^T b. \quad (5.29)$$

因此, 只要求得 $\mathcal{R}(A)$ 的一组标准正交基, 就可以通过式 (5.29) 和式 (5.28) 求得最小二乘问题 (5.1) 的任一解.

现由于 $\text{rank}(A) = r < n$, 2.2 节介绍的 QR 分解式 (2.9) 一般不能产生 $\mathcal{R}(A)$ 的一组标准正交基. 但如果先对 A 的列进行适当的排列使其前 r 列线性无关, 然后再进行 QR 分解, 则仍然可以产生 $\mathcal{R}(A)$ 的一组标准正交基.

设有分解式

$$AP = Q \begin{bmatrix} R_{11} & R_{12} \\ O & O \end{bmatrix} \begin{matrix} r \\ m-r \end{matrix}, \quad (5.30)$$

式中: P 为置换矩阵; Q 为正交矩阵; R_{11} 为非奇异的上三角矩阵. 则 Q 的前 r 列就是 $\mathcal{R}(A)$ 的一组标准正交基.

一旦求出式 (5.30) 的分解式, 则由式 (5.26), 有

$$(Q^T AP)(P^T x) = Q^T b_1,$$

由式 (5.30) 并令 $P^T x = \begin{bmatrix} w \\ z \end{bmatrix}$, 得

$$\begin{bmatrix} R_{11} & R_{12} \\ O & O \end{bmatrix} \begin{bmatrix} w \\ z \end{bmatrix} = Q^T b_1 = \begin{bmatrix} Q_1^T \\ Q_2^T \end{bmatrix} b_1 = \begin{bmatrix} h \\ g \end{bmatrix},$$

即

$$\begin{bmatrix} R_{11}w + R_{12}z \\ 0 \end{bmatrix} = \begin{bmatrix} h \\ g \end{bmatrix}.$$

由此, 得

$$w = R_{11}^{-1}(h - R_{12}z), \quad g = 0.$$

故

$$x = P \begin{bmatrix} w \\ z \end{bmatrix} = P \begin{bmatrix} R_{11}^{-1}(h - R_{12}z) \\ z \end{bmatrix}, \quad z \in \mathbb{R}^{n-r},$$

这就是最小二乘问题 (5.1) 的通解. 注意到上式还可以写为

$$x = x_b + P \begin{bmatrix} -R_{11}^{-1}R_{12} \\ I_{n-r} \end{bmatrix} z, \quad x_b = P \begin{bmatrix} R_{11}^{-1}h \\ 0 \end{bmatrix}, \quad z \in \mathbb{R}^{n-r}, \quad (5.31)$$

式中: x_b 为最小二乘问题的基本解.

下面讨论分解式 (5.30) 的具体实施过程. 类似于分解式 (2.9) 的计算过程, 这一分解可用 Householder 变换和适当的初等列变换相结合逐步求得. 假设对某一正整数 k , 已经求得 $k-1$ 个 Householder 变换 H_1, H_2, \dots, H_{k-1} 和 $k-1$ 个初等变换矩阵 P_1, P_2, \dots, P_{k-1} , 使得

$$R_{k-1} = (H_{k-1} \cdots H_2 H_1) A (P_1 P_2 \cdots P_{k-1})$$

$$= \begin{bmatrix} R_{11}^{(k-1)} & R_{12}^{(k-1)} \\ O & R_{22}^{(k-1)} \end{bmatrix} \begin{matrix} k-1 \\ m-k+1 \\ k-1 & n-k+1 \end{matrix}, \quad (5.32)$$

式中: $R_{11}^{(k-1)}$ 为非奇异的上三角矩阵. 现记

$$R_{22}^{(k-1)} = [v_k^{(k-1)}, v_{k+1}^{(k-1)}, \dots, v_n^{(k-1)}],$$

即 $v_i^{(k-1)}$ 表示 $R_{22}^{(k-1)}$ 的第 $i-k+1$ 列. 下一步是首先确定指标 p ($k \leq p \leq n$) 满足

$$\|v_p^{(k-1)}\|_2 = \max \{ \|v_k^{(k-1)}\|_2, \|v_{k+1}^{(k-1)}\|_2, \dots, \|v_n^{(k-1)}\|_2 \}. \quad (5.33)$$

若 $\|v_p^{(k-1)}\|_2 = 0$, 停止计算. 否则取 P_k 为第 k 列与第 p 列交换的初等变换矩阵, 并确定一个 Householder 变换 $\widetilde{H}_k \in \mathbb{R}^{(m-k+1) \times (m-k+1)}$, 使得

$$\widetilde{H}_k v_p^{(k-1)} = \gamma_{kk} e_1.$$

令 $H_k = \text{diag}(I_{k-1}, \widetilde{H}_k)$, 则有

$$R_k = \begin{bmatrix} R_{11}^{(k)} & R_{12}^{(k)} \\ O & R_{22}^{(k)} \end{bmatrix} \begin{matrix} k \\ m-k \\ k & n-k \end{matrix}, \quad (5.34)$$

式中: $R_{11}^{(k)}$ 为 $k \times k$ 阶非奇异上三角矩阵.

这样, 从 $k=1$ 出发, 依次进行 $r (= \text{rank}(\mathbf{A}))$ 次“列选主元”的 Householder 变换, 即可得到分解式 (5.30). 值得一提的是, 如果只是利用这一分解来求秩亏的最小二乘问题的解 (通常是用来求基本解), 那么正交矩阵 $Q = H_1 H_2 \cdots H_r$ 没有必要显式地计算, 只需将计算过程中每一步的 Householder 变换同步地作用在向量 b 上即可.

此外, 为了减少“列选主元”的计算量, 不必每一步都按照 2-范数的定义去计算式 (5.33) 中的范数, 因为对于任何正交矩阵 $Q \in \mathbb{R}^{l \times l}$ 均成立

$$Q^T x = \begin{bmatrix} \alpha \\ z \end{bmatrix} \begin{matrix} 1 \\ l-1 \end{matrix} \implies \|z\|_2^2 = \|x\|_2^2 - \alpha^2.$$

这样, 可以通过修正旧的范数来得到新的范数, 即

$$\|x^{(j)}\|_2^2 = \|x^{(j-1)}\|_2^2 - a_{kj}^2.$$

综上所述, 可写出计算分解式 (5.30) 的算法并用以求解秩亏最小二乘问题的详细步骤如下.

算法 5.3 (秩亏最小二乘问题的列主元 QR 分解法)

步 1, 输入矩阵 $A = (a_{ij}) \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$. 计算

$$p(j) := j, \quad \sigma_j = A(:, j)^T A(:, j), \quad j = 1, 2, \dots, n.$$

置 $k := 1$.

步 2, 确定 r 使满足 $\sigma_r = \max_{k \leq j \leq n} \sigma_j$. 若 $\sigma_r = 0$, 停算.

步 3, 交换 $p(r)$ 与 $p(k)$, σ_r 与 σ_k 以及 $A(:, r)$ 与 $A(:, k)$.

步 4, 确定一个 Householder 矩阵 \widetilde{H}_k 使得

$$\widetilde{H}_k A(k:m, k) = \gamma_{kk} e_1,$$

并计算

$$[A, b] := \text{diag}(I_{k-1}, \widetilde{H}_k)[A, b], \quad \sigma_j := \sigma_j - a_{kj}^2, \quad j = k+1, \dots, n.$$

步 5, 置 $k := k+1$, 转步 2.

算法 5.3 的运算量是 $2mnr - r^2(m+n) + 2r^3/3$ 次乘除法, 其中 $r = \text{rank}(A)$. 分解式 (5.30) 中的上三角矩阵 R 存储在 A 的上三角部分, 初等变换矩阵 P 由整数向量 $p(1:n)$ 来记录, 即 $P = P_1 P_2 \cdots P_r$, 其中 P_k 是通过单位矩阵交换第 k 列与第 $p(k)$ 列而得到. 而第 k 个 Householder 向量的 $k+1:n$ 分量存储在 $A(k+1:m, k)$ 中.

根据算法 5.3 编制 MATLAB 程序如下:

```
function [x,fval,P]=piv_house_qr(A,b)
%秩亏最小二乘问题的列主元QR分解法
[m,n]=size(A); %Q=eye(m);
A1=A; b1=b;
P=eye(n);
for j=1:n
    c(j)=A(1:m,j)'\*A(1:m,j);
end
[cr,r]=max(c);
for k=1:n
    if (cr<=0), break; end
    c([k r])=c([r k]);
    P(:, [k r])=P(:, [r k]);
    A(1:m, [k r])= A(1:m, [r k]);
    [v,beta]=house(A(k:m,k));
    H=eye(m-k+1)-beta*v*v';
    A(k:m,k:n)=H*A(k:m,k:n);
    b(k:m)=H*b(k:m);
    for j=k+1:n
        c(j)=c(j)-A(k,j)^2;
    end
    [cr,r]=max(c(k+1:n));
```

```

    r=r+k;
end
for i=1:m
    if sum(abs(A(i,:)))<1.0e-10
        rank=i-1; break;
    end
end
R11=A(1:rank,1:rank); c=b(1:rank);
x=P*[R11\c; zeros(rank,1)];
fval=norm(A1*x-b1);

```

例 5.2 用 MATLAB 程序 piv_house_qr.m 求解超定方程组 $Ax = b$ 的最小二乘解, 其中

$$A = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1 & 4 & 5 & 6 \\ 1 & 5 & 6 & 7 \\ 1 & 8 & 9 & 10 \\ 1 & 11 & 12 & 13 \end{bmatrix}, \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}, \quad b = \begin{bmatrix} 11 \\ 13 \\ 15 \\ 18 \\ 20 \end{bmatrix}.$$

解 在 MATLAB 命令窗口输入:

```

>> A=[1 2 3 4; 1 4 5 6; 1 5 6 7; 1 8 9 10; 1 11 12 13];
>> b=[11 13 15 18 20]';
>> [x,fval,P]=piv_house_qr(A,b)

```

即可得计算结果.

此外, 如果希望求出秩亏最小二乘问题的极小范数解, 则还需将分解式 (5.30) 中的 R_{12} 约化为零矩阵. 这可通过 r 次 Householder 变换来完成. 即确定 r 个 Householder 变换 Z_1, \dots, Z_r 和一个置换矩阵 \tilde{P} 使得

$$[R_{11}, R_{12}]Z_1 \cdots Z_r \tilde{P} = [T, O],$$

式中: $T \in \mathbb{R}^{r \times r}$ 为上三角矩阵.

现令 $Z = PZ_1 \cdots Z_r \tilde{P}$, 则有

$$Q^T A Z = \begin{bmatrix} T & O \\ O & O \end{bmatrix} \begin{matrix} r \\ m-r \end{matrix}, \quad (5.35)$$

$r \quad n-r$

由此, 得

$$x_{LS} = Z \begin{bmatrix} T^{-1}c \\ 0 \end{bmatrix},$$

式中: c 为由 $Q^T b$ 前 r 个分量构成的 r 维向量.

5.3.2 奇异值分解法

设 $A \in \mathbb{R}^{m \times n}$ ($m > n$) 的奇异值分解为

$$A = U \Sigma V^T, \quad \Sigma = \begin{bmatrix} \Sigma_r & O \\ O & O \end{bmatrix} \begin{matrix} r \\ m-r \\ r & n-r \end{matrix}, \quad (5.36)$$

式中: $U = [u_1, u_2, \dots, u_m]$ 和 $V = [v_1, v_2, \dots, v_n]$ 为正交矩阵; $\Sigma_r = \text{diag}(\sigma_1, \dots, \sigma_r)$, $\sigma_1 \geq \dots \geq \sigma_r > 0$. 则由定理 5.2 和定理 5.3 可知

$$x_{LS} = A^\dagger b = V \begin{bmatrix} \Sigma_r^{-1} & O \\ O & O \end{bmatrix} U^T = \sum_{i=1}^r \frac{u_i^T b}{\sigma_i} v_i. \quad (5.37)$$

因此, 一旦求出分解式 (5.36), 就可由式 (5.37) 容易地求出最小二乘问题 (5.1) 的极小范数解 x_{LS} .

例 5.3 用 MATLAB 系统自带的奇异值分解函数 `svd.m` 求解超定方程组 $Ax = b$, 其中

$$A = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1 & 4 & 5 & 6 \\ 1 & 5 & 6 & 7 \\ 1 & 8 & 9 & 10 \\ 1 & 11 & 12 & 13 \end{bmatrix}, \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}, \quad b = \begin{bmatrix} 11 \\ 13 \\ 15 \\ 18 \\ 20 \end{bmatrix}.$$

解 编制 MATLAB 程序, 并存成文件名 `ex53.m`:

`%例5.3 ex53.m`

`A=[1 2 3 4; 1 4 5 6; 1 5 6 7; 1 8 9 10; 1 11 12 13];`

`b=[11 13 15 18 20]';`

`[m,n]=size(A);`

`[U,S,V]=svd(A);`

`% x=(V*pinv(S)*U')*b;`

`for i=1:n`

`if abs(S(i,i))<1.0e-6`

`r=i-1; break;`

`end`

`end %r=rank(S);`

`x=zeros(n,1);`

`for i=1:r`

`x=x+(U(:,i))*b/S(i,i))*V(:,i);`

`end`

`x`

`fval=norm(A*x-b)`

然后在命令窗口输入 ex53 即可得计算结果.

观察到求出的最小二乘解与例 5.2 不一样. 原因是本例求出的是极小范数解, 而例 5.2 只是一个基解. 但它们对应的最优值是相同的.

5.4 求解最小二乘问题的迭代方法

对于某些大型稀疏的最小二乘问题, 前面介绍的直接解法不一定有效. 此时, 迭代法应该是最好的选择. 本节讨论大型稀疏最小二乘问题

$$\|Ax - b\|_2 = \min_{z \in \mathbb{R}^n} \|Az - b\|_2 \quad (5.38)$$

解的迭代算法, 包括矩阵分裂迭代法和 Krylov 子空间迭代法. 本节假定矩阵 A 是 $m \times n$ ($m \geq n$) 阶的列满秩实矩阵, b 是 m 维实向量, 即 $A \in \mathbb{R}_n^{m \times n}$, $b \in \mathbb{R}^m$. 于是, 最小二乘问题 (5.38) 有唯一的最小二乘解 x_{LS} . 这里所考虑的迭代算法基于最小二乘问题的法方程

$$A^T A x = A^T b, \quad (5.39)$$

或基于其 KKT 方程

$$\begin{bmatrix} I_m & A \\ A^T & O \end{bmatrix} \begin{bmatrix} r \\ x \end{bmatrix} = \begin{bmatrix} b \\ 0 \end{bmatrix}. \quad (5.40)$$

5.4.1 基于法方程的矩阵分裂迭代法

由于 $A \in \mathbb{R}^{m \times n}$ 是列满秩的, 因此, 原则上讲, 任何适用于求解对称正定方程组的迭代方法都可用于求解最小二乘问题的法方程 (5.39), 从而得到相应的求解最小二乘问题 (5.38) 的迭代算法.

为了方便起见, 将 A 按列分块为 $A = [a_1, a_2, \dots, a_n]$, 并记 $d_i = a_i^T a_i$, $i = 1, 2, \dots, n$. $D = \text{diag}(d_1, d_2, \dots, d_n)$. 于是可将矩阵 $A^T A$ 分裂为

$$A^T A = D - L - L^T, \quad (5.41)$$

这里 $-L$ 是其严格下三角部分.

1. Jacobi 迭代法

在分裂 $A^T A = M - N$ 中, 令 $M = D$, $N = L + L^T$, 则得到 Jacobi 迭代格式为

$$x^{(k+1)} = D^{-1} [(D - A^T A)x^{(k)} + A^T b], \quad k = 0, 1, \dots, \quad (5.42)$$

或等价地, 有

$$x^{(k+1)} = x^{(k)} + D^{-1} A^T (b - Ax^{(k)}), \quad k = 0, 1, \dots, \quad (5.43)$$

显然, Jacobi 迭代法的迭代矩阵为 $G_J = I - D^{-1} A^T A$, 且迭代过程不需要显式地计算 $A^T A$.

在实际计算中, 可以对每个分量进行迭代. 令 $r^{(k)} = b - Ax^{(k)}$, 执行下面的算法:

```

 $\mathbf{r}^{(0)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(0)}; k = 0;$ 
while ( $\|\mathbf{r}^{(k)}\|_2 / \|\mathbf{r}^{(0)}\|_2 \leq \varepsilon$ )
    for  $i = 1:n$ 
         $d_i = \mathbf{a}_i^T \mathbf{a}_i;$ 
         $x_i^{(k+1)} = x_i^{(k)} + \mathbf{a}_i^T \mathbf{r}^{(k)} / d_i;$ 
    end
     $\mathbf{r}^{(k+1)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k+1)};$ 
     $k = k + 1;$ 
end

```

Jacobi 迭代法的 MATLAB 程序如下:

```

function [x,fval,k,time]=nls_jacobi(A,b,x,tol,max_it)
%Jacobi算法求解法方程 $\mathbf{A}'\mathbf{A}\mathbf{x}=\mathbf{A}'\mathbf{b}$ 
tic; n=size(A,2); r=b-A*x;
d=zeros(n,1); nr=norm(A'*r); k=0;
for i=1:n
    d(i)=A(:,i)'\*A(:,i);
end
while (k<=max_it)
    k=k+1;
    for i=1:n
        x(i)=x(i)+A(:,i)'\*r/d(i);
    end
    r=b-A*x; s=A'*r;
    if (norm(s)/nr<tol)
        break;
    end
end
fval=norm(r); time=toc;

```

例 5.4 利用 Jacobi 迭代法的 MATLAB 程序, 计算最小二乘问题 $\min \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2$ 的极小解, 其中

$$\mathbf{A} = \begin{bmatrix} 4 & -1 & 0 & 0 \\ 1 & 4 & -1 & 0 \\ 0 & 1 & 4 & -1 \\ 0 & 0 & 1 & 4 \\ 1 & 3 & 2 & 1 \\ 0 & 2 & 0 & 3 \\ 8 & 2 & 3 & 1 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 9 \\ 12 \\ 11 \\ 13 \\ 17 \\ 15 \\ 19 \end{bmatrix}.$$

取初始点为零向量, 容许误差为 10^{-10} .

解 编写 MATLAB 脚本文件 ex54.m, 并在命令窗口运行之, 迭代 708 次, 得到极小解和极小值

$$\mathbf{x}^* = (1.1019, 2.7905, 1.9907, 2.6509)^T, \quad \|\mathbf{b} - \mathbf{A}\mathbf{x}^*\|_2 = 9.2279.$$

2. Gauss-Seidel 迭代法

在分裂 $\mathbf{A}^T \mathbf{A} = \mathbf{M} - \mathbf{N}$ 中, 令 $\mathbf{M} = \mathbf{D} - \mathbf{L}$, $\mathbf{N} = \mathbf{L}^T$, 则得到 Gauss-Seidel 迭代格式为

$$\mathbf{x}^{(k+1)} = \mathbf{D}^{-1}[\mathbf{L}\mathbf{x}^{(k+1)} + (\mathbf{D} - \mathbf{L} - \mathbf{A}^T \mathbf{A})\mathbf{x}^{(k)} + \mathbf{A}^T \mathbf{b}], \quad k = 0, 1, \dots, \quad (5.44)$$

或等价地, 有

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + (\mathbf{D} - \mathbf{L})^{-1} \mathbf{A}^T (\mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}), \quad k = 0, 1, \dots. \quad (5.45)$$

不难发现, Gauss-Seidel 迭代法的迭代矩阵为 $\mathbf{G}_S = (\mathbf{D} - \mathbf{L})^{-1}(\mathbf{D} - \mathbf{L} - \mathbf{A}^T \mathbf{A}) = (\mathbf{D} - \mathbf{L})^{-1} \mathbf{L}^T$, 并需要显式地计算 $\mathbf{A}^T \mathbf{A}$.

为了避免 $\mathbf{A}^T \mathbf{A}$ 的显式计算, 可在第 k 步引入辅助变量 $\mathbf{z}_1 = \mathbf{x}^{(k)}$, $\mathbf{r}_1 = \mathbf{r}^{(k)}$, 则式 (5.44) 可化为

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \mathbf{D}^{-1}[\mathbf{L}\mathbf{x}^{(k+1)} + \mathbf{A}^T \mathbf{b} + (\mathbf{D} - \mathbf{L} - \mathbf{A}^T \mathbf{A})\mathbf{x}^{(k)}] \\ &= \mathbf{x}^{(k)} + \mathbf{D}^{-1}[\mathbf{L}(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) + \mathbf{A}^T (\mathbf{b} - \mathbf{A}\mathbf{x}^{(k)})] \\ &= \mathbf{z}_1 + \mathbf{D}^{-1}[\mathbf{L}(\mathbf{x}^{(k+1)} - \mathbf{z}_1) + \mathbf{A}^T \mathbf{r}_1]. \end{aligned} \quad (5.46)$$

详细的算法步骤如下:

```

 $\mathbf{r}^{(0)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(0)}; k = 0;$ 
while ( $\|\mathbf{r}^{(k)}\|_2 / \|\mathbf{r}^{(0)}\|_2 \leq \varepsilon$ )
     $\mathbf{r}_1 = \mathbf{r}^{(k)}; \mathbf{z}_1 = \mathbf{x}^{(k)};$ 
    for  $i = 1 : n$ 
         $d_i = \mathbf{a}_i^T \mathbf{a}_i; \alpha_i = \mathbf{a}_i^T \mathbf{r}_i / d_i;$ 
         $\mathbf{z}_{i+1} = \mathbf{z}_i + \alpha_i \mathbf{e}^{(i)};$ 
         $\mathbf{r}_{i+1} = \mathbf{r}_i - \alpha_i \mathbf{a}_i;$ 
    end
     $\mathbf{x}^{(k+1)} = \mathbf{z}_{n+1};$ 
     $\mathbf{r}^{(k+1)} = \mathbf{r}_{n+1};$ 
     $k = k + 1;$ 
end

```

以上算法不再需要显式地计算 $\mathbf{A}^T \mathbf{A}$.

Gauss-Seidel 迭代法的 MATLAB 程序如下:

```

function [x,fval,k,time]=nls_seidel(A,b,x,tol,max_it)
%Gauss--Seidel算法求解方程 $\mathbf{A}'\mathbf{A}\mathbf{x}=\mathbf{A}'\mathbf{b}$ 

```

```

tic;n=size(A,2);
r=b-A*x;d=zeros(n,1);
nr=norm(A'*r);k=0;E=eye(n);
for i=1:n
    d(i)=A(:,i)'\*A(:,i);
end
D=diag(d); Z=[x]; R=[r];
while (k<=max_it)
    k=k+1;
    for i=1:n
        delta=A(:,i)'\*R(:,i)/d(i);
        z=Z(:,i)+delta*E(:,i);
        r=R(:,i)-delta*A(:,i);
        Z=[Z,z]; R=[R,r];
    end
    x=Z(:,n+1);r=b-A*x;
    Z=[x];R=[r];
    if(norm(A'*r)/nr<tol),break;end
end
fval=norm(r); time=toc;

```

例 5.5 利用 Gauss-Seidel 迭代法的 MATLAB 程序, 计算例 5.4 中的最小二乘问题 $\min \|b - Ax\|_2$ 的极小解. 取初始点为零向量, 容许误差为 10^{-10} .

解 编写 MATLAB 脚本文件 ex55.m, 并在命令窗口运行之, 迭代 19 次, 得到极小解和极小值

$$x^* = (1.1019, 2.7905, 1.9907, 2.6509)^T, \quad \|b - Ax^*\|_2 = 9.2279.$$

3. SOR 迭代法

在分裂 $A^T A = M - N$ 中, 令

$$M = \frac{1}{\omega}(D - \omega L), \quad N = \frac{1}{\omega}[(1 - \omega)D + \omega L^T], \quad \omega \neq 0,$$

则得到 SOR 迭代格式为

$$x^{(k+1)} = \omega D^{-1}(Lx^{(k+1)} + L^T x^{(k)} + A^T b) + (1 - \omega)x^{(k)}, \quad k = 0, 1, \dots, \quad (5.47)$$

式中: ω 为松弛因子. 或等价地, 有

$$x^{(k+1)} = x^{(k)} + \omega D^{-1}[L(x^{(k+1)} - x^{(k)}) + A^T(b - Ax^{(k)})], \quad k = 0, 1, \dots. \quad (5.48)$$

上述迭代法也需要显式地计算 $A^T A$, 其迭代矩阵为 $G_\omega = (D - \omega L)^{-1}[(1 - \omega)D + \omega L^T]$.

同样, 为了避免显式计算 $A^T A$, 可在第 k 步引入辅助变量 $z_1 = x^{(k)}$, $r_1 = r^{(k)}$, 则式 (5.47) 可化为

$$\begin{aligned} x^{(k+1)} &= x^{(k)} + \omega D^{-1} [L(x^{(k+1)} - x^{(k)}) + A^T(b - Ax^{(k)})] \\ &= z_1 + \omega D^{-1} [L(x^{(k+1)} - z_1) + A^T r_1]. \end{aligned} \quad (5.49)$$

详细的算法步骤如下:

```

 $r^{(0)} = b - Ax^{(0)}$ ;  $k = 0$ ;
while ( $\|r^{(k)}\|_2 / \|r^{(0)}\|_2 \leq \varepsilon$ )
     $r_1 = r^{(k)}$ ;  $z_1 = x^{(k)}$ ;
    for  $i = 1 : n$ 
         $d_i = a_i^T a_i$ ;  $\alpha_i = \omega a_i^T r_i / d_i$ ;
         $z_{i+1} = z_i + \alpha_i e^{(i)}$ ;
         $r_{i+1} = r_i - \alpha_i a_i$ ;
    end
     $x^{(k+1)} = z_{n+1}$ ,  $r^{(k+1)} = r_{n+1}$ ;
     $k = k + 1$ ;
end

```

注 5.4 关于上述三种迭代法的收敛性定理, 跟第 3 章中求解线性方程组的 Jacobi 迭代法、Gauss-Seidel 迭代法和 SOR 迭代法相类似, 只需将此处的 $A^T A$ 和 $A^T b$ 分别视作方程 $Ax = b$ 中的 A 和 b 即可得到类似的结论。

SOR 迭代法的 MATLAB 程序如下:

```

function [x,fval,k,time]=nls_sor(A,b,w,x,tol,max_it)
%SOR迭代法求解法方程A'Ax=A'b
tic;n=size(A,2);r=b-A*x;d=zeros(n,1);
nr=norm(A'*r);k=0; E=eye(n);
for(i=1:n),d(i)=A(:,i)'*A(:,i);end
D=diag(d);Z=[x];R=[r];
while (k<=max_it)
    k=k+1;
    for i=1:n
        delta=w*A(:,i)'*R(:,i)/d(i);
        z=Z(:,i)+delta*E(:,i);
        r=R(:,i)-delta*A(:,i);
        Z=[Z,z]; R=[R,r];
    end
    x=Z(:,n+1);r=b-A*x;Z=[x];R=[r];
    if(norm(A'*r)/nr<tol),break;end
end
fval=norm(r);time=toc;

```

例 5.6 利用 SOR 迭代法的 MATLAB 程序, 计算最小二乘问题 $\min \|b - Ax\|_2$ 的极小解, 其中

$$A = \begin{bmatrix} 1 & 23.73 & 5.49 & 1.21 \\ 1 & 22.34 & 4.32 & 1.35 \\ 1 & 28.84 & 5.04 & 1.92 \\ 1 & 27.67 & 4.72 & 1.49 \\ 1 & 20.83 & 5.35 & 1.56 \\ 1 & 22.27 & 4.27 & 1.50 \\ 1 & 27.57 & 5.25 & 1.85 \\ 1 & 28.01 & 4.62 & 1.51 \end{bmatrix}, \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}, \quad b = \begin{bmatrix} 15.02 \\ 12.62 \\ 14.86 \\ 13.98 \\ 15.91 \\ 12.47 \\ 15.80 \\ 14.32 \end{bmatrix}.$$

取初始点为零向量, 松弛因子 $\omega = 1.06$, 容许误差为 10^{-6} .

解 编写 MATLAB 脚本文件 ex56.m, 并在命令窗口运行之, 迭代 1134 次, 得到极小解和极小值

$$x^* = (-0.0054, 0.0169, 2.4468, 1.2954)^T, \quad \|b - Ax^*\|_2 = 0.9959.$$

5.4.2 基于法方程的共轭梯度法

由于 $A \in \mathbb{R}^{m \times n}$ ($m \geq n$) 列满秩, 故最小二乘问题的法方程 (5.39) 是对称正定方程组. 因此, 法方程 (5.39) 等价于极小化问题

$$\min f(x) = \frac{1}{2} x^T (A^T A) x - (A^T b)^T x. \quad (5.50)$$

事实上, $f(x)$ 的梯度为

$$\nabla f(x) = A^T A x - A^T b,$$

且对任意给定的非零向量 p 和实数 α , 有

$$f(x + \alpha p) = f(x) + \alpha \nabla f(x)^T p + \frac{1}{2} \alpha^2 p^T A^T A p.$$

若 x^* 是法方程 (5.39) 的解, 则有 $\nabla f(x^*) = 0$. 因此, 对任意非零向量 $p \in \mathbb{R}^n$, 有

$$f(x^* + \alpha p) \begin{cases} > f(x^*), & \alpha \neq 0, \\ = f(x^*), & \alpha = 0. \end{cases}$$

故 x^* 是 $f(x)$ 的极小点. 反之, 因 $A^T A$ 对称正定, 可知 $f(x)$ 在 \mathbb{R}^n 中有唯一的极小点. 若 x^* 是 $f(x)$ 的极小点, 则必有 $\nabla f(x^*) = A^T A x^* - A^T b = 0$, 即 x^* 是法方程 (5.39) 的解.

类似于用共轭梯度法求解对称正定方程组 $Ax = b$ 的思想, 设 $x^{(0)} \in \mathbb{R}^n$ 是任意给定的一个初始向量, 对于 $k = 0, 1, 2, \dots$, 从 $x^{(k)}$ 出发, 沿方向 $p^{(k)}$ 求函数 $f(x)$ 在直线 $x = x^{(k)} + \alpha p^{(k)}$ 上的极小点, 得

$$x^{(k+1)} = x^{(k)} + \alpha_k p^{(k)}, \quad r^{(k)} = b - Ax^{(k)}, \quad (5.51)$$

$$s^{(k)} = A^T r^{(k)}, \quad \alpha_k = \frac{(s^{(k)}, p^{(k)})}{(p^{(k)}, A^T A p^{(k)})}, \quad (5.52)$$

称向量 $p^{(k)}$ 为搜索方向. 若向量组 $p^{(0)}, p^{(1)}, \dots, p^{(k-1)}$ 满足

$$(p^{(i)}, A^T A p^{(j)}) = 0, \quad i \neq j, \quad (5.53)$$

且 $p^{(k)} \neq 0, k = 0, 1, \dots, n-1$, 则称向量组 $\{p^{(k)}\}$ 为 \mathbb{R}^n 中关于 $A^T A$ 的一个共轭向量组, 称迭代法 (5.51) 为共轭方向法. 特别地, 若取 $p^{(0)} = s^{(0)}$,

$$p^{(k+1)} = s^{(k+1)} + \beta_k p^{(k)}, \quad \beta_k = -\frac{(s^{(k+1)}, A^T A p^{(k)})}{(p^{(k)}, A^T A p^{(k)})}, \quad (5.54)$$

则称为共轭梯度法.

由式 (5.51)、式 (5.52) 及式 (5.54) 可知, 若存在 $k \geq 0$, 使得 $s^{(k)} = 0$, 则 $x^{(k)}$ 为最小二乘问题的解, 且有 $\alpha_k = \beta_k = 0, s^{(k+1)} = p^{(k+1)} = 0$. 此外, 在上述的共轭梯度法中每次迭代需要计算 4 次矩阵与向量的乘法, 因此有必要对 α_k 和 β_k 的计算公式进行化简.

定理 5.6 在共轭梯度法 (5.51)、(5.52) 和 (5.54) 中, 若 $k > 0, s^{(k)} \neq 0$, 则有

$$\begin{cases} (s^{(k)}, s^{(i)}) = (s^{(k)}, p^{(i)}) = (p^{(k)}, A^T A p^{(i)}) = 0, & 0 \leq i < k, \\ (p^{(k)}, s^{(i)}) = (s^{(k)}, s^{(k)}), & 0 \leq i \leq k. \end{cases} \quad (5.55)$$

证明 由定理的条件, 有 $\alpha_i \neq 0, i = 0, 1, \dots, k$. 根据迭代公式, 得

$$\begin{cases} s^{(k+1)} = A^T r^{(k+1)} = A^T (b - Ax^{(k+1)}) = s^{(k)} - \alpha_k A^T A p^{(k)}, \\ p^{(k+1)} = s^{(k+1)} + \beta_k p^{(k)} = s^{(k)} - \alpha_k A^T A p^{(k)} + \beta_k p^{(k)}. \end{cases} \quad (5.56)$$

下面用归纳法证明定理的结论. 注意到

$$\begin{aligned} (s^{(1)}, s^{(0)}) &= (s^{(1)}, p^{(0)}) = (s^{(0)}, p^{(0)}) - \alpha_0 (A^T A p^{(0)}, p^{(0)}) = 0, \\ (p^{(1)}, A^T A p^{(0)}) &= (s^{(1)}, A^T A p^{(0)}) + \beta_0 (p^{(0)}, A^T A p^{(0)}) = 0, \\ (p^{(1)}, s^{(0)}) &= (p^{(1)}, s^{(1)} + \alpha_0 A^T A p^{(0)}) = (p^{(1)}, s^{(1)}) \\ &= (s^{(1)} + \beta_0 p^{(0)}, s^{(1)}) = (s^{(1)}, s^{(1)}), \end{aligned}$$

即 $k = 1$ 时结论成立. 现假定一直到 $k (> 1)$ 时, 结论成立. 则

$$(s^{(k+1)}, p^{(k)}) = (s^{(k)}, p^{(k)}) - \alpha_k (A^T A p^{(k)}, p^{(k)}) = 0,$$

$$\begin{aligned}
(\mathbf{p}^{(k+1)}, \mathbf{A}^T \mathbf{A} \mathbf{p}^{(k)}) &= (\mathbf{s}^{(k+1)}, \mathbf{A}^T \mathbf{A} \mathbf{p}^{(k)}) + \beta_k (\mathbf{p}^{(k)}, \mathbf{A}^T \mathbf{A} \mathbf{p}^{(k)}) = 0, \\
(\mathbf{s}^{(k+1)}, \mathbf{s}^{(k)}) &= (\mathbf{s}^{(k+1)}, \mathbf{p}^{(k)} - \beta_{k-1} \mathbf{p}^{(k-1)}) = -\beta_{k-1} (\mathbf{s}^{(k+1)}, \mathbf{p}^{(k-1)}) \\
&= -\beta_{k-1} (\mathbf{s}^{(k)} - \alpha_k \mathbf{A}^T \mathbf{A} \mathbf{p}^{(k)}, \mathbf{p}^{(k-1)}) = 0.
\end{aligned}$$

当 $i < k$ 时, 有

$$\begin{aligned}
(\mathbf{s}^{(k+1)}, \mathbf{p}^{(i)}) &= (\mathbf{s}^{(k)}, \mathbf{p}^{(i)}) - \alpha_k (\mathbf{A}^T \mathbf{A} \mathbf{p}^{(k)}, \mathbf{p}^{(i)}) = 0, \\
(\mathbf{p}^{(k+1)}, \mathbf{A}^T \mathbf{A} \mathbf{p}^{(i)}) &= (\mathbf{s}^{(k+1)}, \mathbf{A}^T \mathbf{A} \mathbf{p}^{(i)}) + \beta_k (\mathbf{p}^{(k)}, \mathbf{A}^T \mathbf{A} \mathbf{p}^{(i)}) \\
&= (\mathbf{s}^{(k+1)}, \mathbf{A}^T \mathbf{A} \mathbf{p}^{(i)}) = (\mathbf{s}^{(k+1)}, (\mathbf{s}^{(i)} - \mathbf{s}^{(i+1)})/\alpha_i) = 0, \\
(\mathbf{s}^{(k+1)}, \mathbf{s}^{(i)}) &= (\mathbf{s}^{(k+1)}, \mathbf{p}^{(i)} - \beta_{i-1} \mathbf{p}^{(i-1)}) = -\beta_{i-1} (\mathbf{s}^{(k+1)}, \mathbf{p}^{(i-1)}) \\
&= -\beta_{i-1} (\mathbf{s}^{(k)} - \alpha_k \mathbf{A}^T \mathbf{A} \mathbf{p}^{(k)}, \mathbf{p}^{(i-1)}) = 0.
\end{aligned}$$

又当 $i \leq k+1$ 时, 有

$$\begin{aligned}
(\mathbf{p}^{(k+1)}, \mathbf{s}^{(i)}) &= \left(\mathbf{p}^{(k+1)}, \mathbf{s}^{(k+1)} + \sum_{l=i}^k \alpha_l \mathbf{A}^T \mathbf{A} \mathbf{p}^{(l)} \right) = (\mathbf{p}^{(k+1)}, \mathbf{s}^{(k+1)}) \\
&= (\mathbf{s}^{(k+1)} + \beta_k \mathbf{p}^{(k)}, \mathbf{s}^{(k+1)}) = (\mathbf{s}^{(k+1)}, \mathbf{s}^{(k+1)}).
\end{aligned}$$

由归纳法原理, 定理得证. □

利用式 (5.56) 的第 1 式, 当 $\alpha_k \neq 0$ 时, 有

$$\mathbf{A}^T \mathbf{A} \mathbf{p}^{(k)} = \frac{\mathbf{s}^{(k)} - \mathbf{s}^{(k+1)}}{\alpha_k},$$

代入 α_k 和 β_k 的表达式并化简, 得

$$\alpha_k = \frac{\|\mathbf{s}^{(k)}\|_2^2}{\|\mathbf{A} \mathbf{p}^{(k)}\|_2^2}, \quad \beta_k = \frac{\|\mathbf{s}^{(k+1)}\|_2^2}{\|\mathbf{s}^{(k)}\|_2^2}. \quad (5.57)$$

算法 5.4 (CGLS) 给定矩阵 $\mathbf{A} \in \mathbb{R}^{m \times n}$, 向量 $\mathbf{b} \in \mathbb{R}^m$ 和容许误差 $\varepsilon > 0$. 本算法计算向量 $\mathbf{x}^{(k)}$, 使得 $\|\mathbf{s}^{(k)}\|_2 / \|\mathbf{s}^{(0)}\|_2 \leq \varepsilon$, 其中 $\mathbf{r}^{(k)} = \mathbf{b} - \mathbf{A} \mathbf{x}^{(k)}$, $\mathbf{s}^{(k)} = \mathbf{A}^T \mathbf{r}^{(k)}$.

取初始向量 $\mathbf{x}^{(0)} \in \mathbb{R}^n$;

计算 $\mathbf{r}^{(0)} = \mathbf{b} - \mathbf{A} \mathbf{x}^{(0)}$; $\mathbf{p}^{(0)} = \mathbf{s}^{(0)} = \mathbf{A}^T \mathbf{r}^{(0)}$; $\gamma_0 = \|\mathbf{s}^{(0)}\|_2^2$; $k = 0$;

while ($\|\mathbf{s}^{(k)}\|_2 / \|\mathbf{s}^{(0)}\|_2 \leq \varepsilon$)

$\mathbf{q}^{(k)} = \mathbf{A} \mathbf{p}^{(k)}$; $\alpha_k = \gamma_k / \|\mathbf{q}^{(k)}\|_2^2$;

$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}$;

$\mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} - \alpha_k \mathbf{q}^{(k)}$; $\mathbf{s}^{(k+1)} = \mathbf{A}^T \mathbf{r}^{(k+1)}$;

$\gamma_{k+1} = \|\mathbf{s}^{(k+1)}\|_2^2$; $\beta_k = \gamma_{k+1} / \gamma_k$;

$\mathbf{p}^{(k+1)} = \mathbf{s}^{(k+1)} + \beta_k \mathbf{p}^{(k)}$;

$k = k + 1$;

end

算法 5.4 的有限终止性、收敛性和收敛速度的分析与求解对称正定方程组 $Ax = b$ 的共轭梯度法相类似, 此处不再赘述. 此外, 不难发现算法 5.4 的每次迭代已经减少到只需计算两次矩阵与向量的乘法.

算法 5.4 的 MATLAB 程序如下:

```
function [x,fval,k,time]=nls_cg(A,b,x,tol,max_it)
%CG算法求解法方程A'Ax=A'b
tic; n=size(A,2); r=b-A*x;
s=A'*r; p=s; gama=s'*s;
nr=norm(s); k=0;
while (k<=max_it)
    k=k+1; q=A*p; alpha=gama/(q'*q);
    x=x+alpha*p; r=r-alpha*q; s=A'*r;
    if (norm(s)/nr<tol), break; end
    gama1=s'*s; beta=gama1/gama;
    p=s+beta*p; gama=gama1;
end
fval=norm(r); time=toc;
```

例 5.7 利用 CG 迭代法的 MATLAB 程序, 计算例 5.6 中的最小二乘问题 $\min \|b - Ax\|_2$ 的极小解, 并与 SOR 迭代法进行比较 (取初始点为零向量, 松弛因子 $\omega = 1.06$, 容许误差 $\varepsilon = 10^{-10}$).

解 编写 MATLAB 脚本文件 ex57.m, 并在命令窗口运行之, 得到数值结果如表 5.1 所示.

表 5.1 CG 法与 SOR 法的数值结果

算法	迭代次数	CPU时间	极小解	极小值
SOR	2889	0.2275	$(-0.0309, 0.0171, 2.4509, 1.2954)^T$	0.9959
CG	5	0.0031	$(-0.0309, 0.0171, 2.4509, 1.2954)^T$	0.9959

5.4.3 基于 KKT 方程的 SOR 类迭代法

由于 SOR 迭代法要求方程组系数矩阵的对角元均不为零, 故不能直接将 SOR 迭代法用于求解最小二乘问题 (5.38) 的 KKT 方程 (5.40), 但可以通过适当的等价变形使之满足 SOR 迭代法的条件.

事实上, KKT 方程 (5.40) 等价于

$$\begin{bmatrix} A & I_m \\ O & A^T \end{bmatrix} \begin{bmatrix} x \\ r \end{bmatrix} = \begin{bmatrix} b \\ 0 \end{bmatrix}. \quad (5.58)$$

由于 $A \in \mathbb{R}_n^{m \times n}$, 不失一般性, 可将 A 进行如下分块:

$$A = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} \begin{matrix} n \\ m-n \end{matrix}, \quad (5.59)$$

式中: $A_1 \in \mathbb{R}^{n \times n}$ 为非奇异阵. 再将 r 和 b 进行相应的分块, 即

$$r = \begin{bmatrix} v \\ w \end{bmatrix}, \quad b = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix},$$

式中: $v, b_1 \in \mathbb{R}^n$, $w, b_2 \in \mathbb{R}^{m-n}$. 则式 (5.58) 变成

$$\begin{bmatrix} A_1 & I_n & O \\ A_2 & O & I_{m-n} \\ O & A_1^T & A_2^T \end{bmatrix} \begin{bmatrix} x \\ v \\ w \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ 0 \end{bmatrix}. \quad (5.60)$$

式 (5.60) 可进一步等价变形为

$$\begin{bmatrix} A_1 & O & I_n \\ A_2 & I_{m-n} & O \\ O & A_2^T & A_1^T \end{bmatrix} \begin{bmatrix} x \\ w \\ v \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ 0 \end{bmatrix}. \quad (5.61)$$

记

$$B = \begin{bmatrix} A_1 & O & I_n \\ A_2 & I_{m-n} & O \\ O & A_2^T & A_1^T \end{bmatrix}, \quad y = \begin{bmatrix} x \\ w \\ v \end{bmatrix}, \quad c = \begin{bmatrix} b_1 \\ b_2 \\ 0 \end{bmatrix},$$

则方程组 (5.61) 可写成

$$By = c. \quad (5.62)$$

由于 A_1 是非奇异的, 不难验证矩阵 B 也是非奇异的.

现按照矩阵 B 的自然分块, 有

$$D = \begin{bmatrix} A_1 & O & O \\ O & I_{m-n} & O \\ O & O & A_1^T \end{bmatrix}, \quad L = - \begin{bmatrix} O & O & O \\ A_2 & O & O \\ O & A_2^T & O \end{bmatrix}, \quad U = - \begin{bmatrix} O & O & I_n \\ O & O & O \\ O & O & O \end{bmatrix}. \quad (5.63)$$

以及

$$\mathcal{L} = D^{-1}L = \begin{bmatrix} O & O & O \\ -A_2 & O & O \\ O & -A_1^{-T}A_2^T & O \end{bmatrix}, \quad \mathcal{U} = D^{-1}U = \begin{bmatrix} O & O & -A_1^{-1} \\ O & O & O \\ O & O & O \end{bmatrix}. \quad (5.64)$$

由此, 可得以下结论:

(1) Jacobi 迭代格式为

$$y^{(k+1)} = \mathcal{G}_J y^{(k)} + f_J, \quad (5.65)$$

式中: 迭代矩阵 \mathcal{G}_J 和右端向量 \mathbf{f}_J 分别为

$$\mathcal{G}_J = \mathcal{L} + \mathcal{U} = \begin{bmatrix} \mathbf{O} & \mathbf{O} & -\mathbf{A}_1^{-1} \\ -\mathbf{A}_2 & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & -(\mathbf{A}_2\mathbf{A}_1^{-1})^T & \mathbf{O} \end{bmatrix}, \quad \mathbf{f}_J = \mathbf{D}^{-1}\mathbf{c} = \begin{bmatrix} \mathbf{A}_1^{-1}\mathbf{b}_1 \\ \mathbf{b}_2 \\ \mathbf{0} \end{bmatrix}. \quad (5.66)$$

详细的迭代格式为

$$\begin{cases} \mathbf{x}^{(k+1)} = \mathbf{A}_1^{-1}(\mathbf{b}_1 - \mathbf{v}^{(k)}), \\ \mathbf{w}^{(k+1)} = \mathbf{b}_2 - \mathbf{A}_2\mathbf{x}^{(k)}, \\ \mathbf{v}^{(k+1)} = -\mathbf{A}_1^{-T}\mathbf{A}_2^T\mathbf{w}^{(k)}, \end{cases} \quad k = 0, 1, 2, \dots \quad (5.67)$$

此外, 若记 $\mathbf{A}_3 = \mathbf{A}_2\mathbf{A}_1^{-1}$, 直接计算, 得

$$\begin{aligned} \mathcal{G}_J^3 &= - \begin{bmatrix} \mathbf{A}_1^{-1}\mathbf{A}_3^T\mathbf{A}_2 & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{A}_2\mathbf{A}_1^{-1}\mathbf{A}_3^T & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{A}_3^T\mathbf{A}_2\mathbf{A}_1^{-1} \end{bmatrix} \\ &= - \begin{bmatrix} \mathbf{A}_1^{-1}\mathbf{A}_3^T\mathbf{A}_3\mathbf{A}_1 & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{A}_3\mathbf{A}_3^T & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{A}_3^T\mathbf{A}_3 \end{bmatrix}. \end{aligned} \quad (5.68)$$

(2) SOR 迭代格式为

$$\mathbf{y}^{(k+1)} = \mathcal{G}_\omega \mathbf{y}^{(k)} + \mathbf{f}_\omega, \quad (5.69)$$

式中: 迭代矩阵 \mathcal{G}_ω 和右端向量 \mathbf{f}_ω 分别为

$$\mathcal{G}_\omega = (\mathbf{I} - \omega\mathcal{L})^{-1}[(1-\omega)\mathbf{I} + \omega\mathcal{U}], \quad \mathbf{f}_\omega = \omega(\mathbf{I} - \omega\mathcal{L})^{-1}\mathbf{f}_J. \quad (5.70)$$

仍记 $\mathbf{A}_3 = \mathbf{A}_2\mathbf{A}_1^{-1}$, 直接计算, 得

$$\begin{aligned} \mathcal{G}_\omega &= \begin{bmatrix} \mathbf{I} & \mathbf{O} & \mathbf{O} \\ \omega\mathbf{A}_2 & \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \omega\mathbf{A}_3^T & \mathbf{I} \end{bmatrix}^{-1} \begin{bmatrix} (1-\omega)\mathbf{I} & \mathbf{O} & -\omega\mathbf{A}_1^{-1} \\ \mathbf{O} & (1-\omega)\mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & (1-\omega)\mathbf{I} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{I} & \mathbf{O} & \mathbf{O} \\ -\omega\mathbf{A}_2 & \mathbf{I} & \mathbf{O} \\ \omega^2\mathbf{A}_3^T\mathbf{A}_2 & -\omega\mathbf{A}_3^T & \mathbf{I} \end{bmatrix} \begin{bmatrix} (1-\omega)\mathbf{I} & \mathbf{O} & -\omega\mathbf{A}_1^{-1} \\ \mathbf{O} & (1-\omega)\mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & (1-\omega)\mathbf{I} \end{bmatrix} \\ &= \begin{bmatrix} (1-\omega)\mathbf{I} & \mathbf{O} & -\omega\mathbf{A}_1^{-1} \\ -\omega(1-\omega)\mathbf{A}_2 & (1-\omega)\mathbf{I} & \omega^2\mathbf{A}_3 \\ \omega(1-\omega)^2\mathbf{A}_3^T\mathbf{A}_2 & -\omega(1-\omega)\mathbf{A}_3^T & (1-\omega)\mathbf{I} - \omega^3\mathbf{A}_3^T\mathbf{A}_3 \end{bmatrix} \\ &:= (1-\omega)\mathbf{I} + \mathcal{K}_\omega, \end{aligned}$$

这里

$$\begin{aligned}
 \mathcal{K}_\omega &= \begin{bmatrix} \mathbf{O} & \mathbf{O} & -\omega \mathbf{A}_1^{-1} \\ -\omega(1-\omega)\mathbf{A}_2 & \mathbf{O} & \omega^2 \mathbf{A}_3 \\ \omega(1-\omega)^2 \mathbf{A}_3^T \mathbf{A}_2 & -\omega(1-\omega)\mathbf{A}_3^T & -\omega^3 \mathbf{A}_3^T \mathbf{A}_3 \end{bmatrix}, \\
 \mathbf{f}_\omega &= \omega \begin{bmatrix} \mathbf{I} & \mathbf{O} & \mathbf{O} \\ \omega \mathbf{A}_2 & \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \omega \mathbf{A}_3^T & \mathbf{I} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{A}_1^{-1} \mathbf{b}_1 \\ \mathbf{b}_2 \\ \mathbf{0} \end{bmatrix} \\
 &= \omega \begin{bmatrix} \mathbf{I} & \mathbf{O} & \mathbf{O} \\ -\omega \mathbf{A}_2 & \mathbf{I} & \mathbf{O} \\ \omega^2 \mathbf{A}_3^T \mathbf{A}_2 & -\omega \mathbf{A}_3^T & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A}_1^{-1} \mathbf{b}_1 \\ \mathbf{b}_2 \\ \mathbf{0} \end{bmatrix} \\
 &= \begin{bmatrix} \omega \mathbf{A}_1^{-1} \mathbf{b}_1 \\ \omega(\mathbf{b}_2 - \omega \mathbf{A}_3 \mathbf{b}_1) \\ -\omega^2 \mathbf{A}_3^T (\mathbf{b}_2 - \omega \mathbf{A}_3 \mathbf{b}_1) \end{bmatrix}.
 \end{aligned}$$

则详细的 SOR 迭代格式为

$$\begin{cases} \mathbf{x}^{(k+1)} = (1-\omega)\mathbf{x}^{(k)} + \omega \mathbf{A}_1^{-1}(\mathbf{b}_1 - \mathbf{v}^{(k)}), \\ \mathbf{w}^{(k+1)} = (1-\omega)\mathbf{w}^{(k)} - \omega(1-\omega)\mathbf{A}_2 \mathbf{x}^{(k)} - \omega^2 \mathbf{A}_3(\mathbf{b}_1 - \mathbf{v}^{(k)}) + \omega \mathbf{b}_2, \\ \mathbf{v}^{(k+1)} = (1-\omega)\mathbf{v}^{(k)} + \omega(1-\omega)^2 \mathbf{A}_3^T \mathbf{A}_2 \mathbf{x}^{(k)} - \omega(1-\omega)\mathbf{A}_3^T \mathbf{w}^{(k)} \\ \quad - \omega^2 \mathbf{A}_3^T [\mathbf{b}_2 - \omega \mathbf{A}_3(\mathbf{b}_1 - \mathbf{v}^{(k)})], \quad k = 0, 1, 2, \dots \end{cases} \quad (5.71)$$

特别地, 在上式中取 $\omega = 1$, 可得 Gauss-Seidel 迭代格式:

$$\begin{cases} \mathbf{x}^{(k+1)} = \mathbf{A}_1^{-1}(\mathbf{b}_1 - \mathbf{v}^{(k)}), \\ \mathbf{w}^{(k+1)} = \mathbf{b}_2 - \mathbf{A}_3(\mathbf{b}_1 - \mathbf{v}^{(k)}), \\ \mathbf{v}^{(k+1)} = -\mathbf{A}_3^T [\mathbf{b}_2 - \mathbf{A}_3(\mathbf{b}_1 - \mathbf{v}^{(k)})], \quad k = 0, 1, 2, \dots \end{cases} \quad (5.72)$$

(3) SSOR 的迭代矩阵 \mathcal{H}_ω 和右端向量 \mathbf{g}_ω 分别为

$$\begin{aligned}
 \mathcal{H}_\omega &= (\mathbf{I} - \omega \mathcal{U})^{-1} [(1-\omega)\mathbf{I} + \omega \mathcal{L}] (\mathbf{I} - \omega \mathcal{L})^{-1} [(1-\omega)\mathbf{I} + \omega \mathcal{U}] \\
 &= (\mathbf{I} - \omega \mathcal{U})^{-1} (\mathbf{I} - \omega \mathcal{L})^{-1} [(1-\omega)\mathbf{I} + \omega \mathcal{L}] [(1-\omega)\mathbf{I} + \omega \mathcal{U}], \quad (5.73)
 \end{aligned}$$

$$\mathbf{g}_\omega = \omega(2-\omega)(\mathbf{I} - \omega \mathcal{U})^{-1} (\mathbf{I} - \omega \mathcal{L})^{-1} \mathbf{f}_J. \quad (5.74)$$

显然, \mathcal{H}_ω 相似于

$$\begin{aligned}
 \hat{\mathcal{H}}_\omega &= (\mathbf{I} - \omega \mathcal{L})^{-1} [(1-\omega)\mathbf{I} + \omega \mathcal{L}] [(1-\omega)\mathbf{I} + \omega \mathcal{U}] (\mathbf{I} - \omega \mathcal{U})^{-1} \\
 &= (\mathbf{I} - \omega \mathcal{L})^{-1} [(1-\omega)\mathbf{I} + \omega \mathcal{L}] (\mathbf{I} - \omega \mathcal{U})^{-1} [(1-\omega)\mathbf{I} + \omega \mathcal{U}] \\
 &:= \mathcal{W}(\mathcal{L})\mathcal{W}(\mathcal{U}), \quad (5.75)
 \end{aligned}$$

这里

$$\mathcal{W}(\mathcal{Z}) = (\mathbf{I} - \omega \mathcal{Z})^{-1}[(1 - \omega)\mathbf{I} + \omega \mathcal{Z}]. \quad (5.76)$$

直接计算, 得

$$\mathcal{W}(\mathcal{L}) = \begin{bmatrix} (1 - \omega)\mathbf{I} & \mathbf{O} & \mathbf{O} \\ -\omega(2 - \omega)\mathbf{A}_2 & (1 - \omega)\mathbf{I} & \mathbf{O} \\ \omega^2(2 - \omega)\mathbf{A}_3^T \mathbf{A}_2 & -\omega(2 - \omega)\mathbf{A}_3^T & (1 - \omega)\mathbf{I} \end{bmatrix},$$

$$\mathcal{W}(\mathcal{U}) = \begin{bmatrix} (1 - \omega)\mathbf{I} & \mathbf{O} & -\omega(2 - \omega)\mathbf{A}_1^{-1} \\ \mathbf{O} & (1 - \omega)\mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & (1 - \omega)\mathbf{I} \end{bmatrix}.$$

于是, 有

$$\widehat{\mathcal{H}}_\omega = \mathcal{W}(\mathcal{L})\mathcal{W}(\mathcal{U}) = (1 - \omega)^2\mathbf{I} + \mathcal{T}_\omega, \quad (5.77)$$

式中:

$$\mathcal{T}_\omega = \begin{bmatrix} \mathbf{O} & \mathbf{O} & -\sigma\mathbf{A}_1^{-1} \\ -\sigma\mathbf{A}_2 & \mathbf{O} & \omega^2(2 - \omega^2)\mathbf{A}_3 \\ \omega\sigma\mathbf{A}_3^T \mathbf{A}_2 & -\sigma\mathbf{A}_3^T & -\omega^2(2 - \omega)^2\mathbf{A}_3^T \mathbf{A}_3 \end{bmatrix},$$

这里 $\sigma = \omega(1 - \omega)(2 - \omega)$.

$$\begin{aligned} \mathbf{g}_\omega &= \omega(2 - \omega)(\mathbf{I} - \omega\mathcal{U})^{-1}(\mathbf{I} - \omega\mathcal{L})^{-1}\mathbf{f}_J \\ &= \omega(2 - \omega) \begin{bmatrix} \mathbf{I} & \mathbf{O} & \omega\mathbf{A}_1^{-1} \\ \mathbf{O} & \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{I} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{I} & \mathbf{O} & \mathbf{O} \\ \omega\mathbf{A}_2 & \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \omega\mathbf{A}_3^T & \mathbf{I} \end{bmatrix}^{-1} \mathbf{f}_J \\ &= \omega(2 - \omega) \begin{bmatrix} \mathbf{I} & \mathbf{O} & -\omega\mathbf{A}_1^{-1} \\ \mathbf{O} & \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{O} & \mathbf{O} \\ -\omega\mathbf{A}_2 & \mathbf{I} & \mathbf{O} \\ \omega^2\mathbf{A}_3^T \mathbf{A}_2 & -\omega\mathbf{A}_3^T & \mathbf{I} \end{bmatrix} \mathbf{f}_J \\ &= \omega(2 - \omega) \begin{bmatrix} \mathbf{A}_1^{-1}\mathbf{b}_1 + \omega^2\mathbf{A}_3^T \mathbf{A}_3\mathbf{b}_1 + \omega^2\mathbf{A}_1^{-1}\mathbf{A}_3^T \mathbf{b}_2 \\ \mathbf{b} - \omega\mathbf{A}_3\mathbf{b}_1 \\ \omega^2\mathbf{A}_3^T \mathbf{A}_3\mathbf{b}_1 - \omega\mathbf{A}_3^T \mathbf{b}_2 \end{bmatrix}. \end{aligned}$$

据此即可写出 SSOR 迭代法的迭代格式

$$\mathbf{y}^{(k+1)} = \mathcal{T}_\omega \mathbf{y}^{(k)} + \mathbf{g}_\omega. \quad (5.78)$$

上述 Jacobi 迭代、SOR 迭代和 SSOR 迭代的收敛性分析可参考文献 [23].

5.4.4 基于 KKT 方程的 HSS 迭代法

现在考虑最小二乘问题 (5.38) 的 KKT 方程

$$\begin{bmatrix} I & A \\ A^T & O \end{bmatrix} \begin{bmatrix} r \\ x \end{bmatrix} = \begin{bmatrix} b \\ 0 \end{bmatrix} \quad (5.79)$$

的迭代解法. 由于式 (5.79) 的系数矩阵是对称不定的, 为了设计更有效的迭代算法, 将其改写成如下的等价形式:

$$Bz := \begin{bmatrix} I & A \\ -A^T & O \end{bmatrix} \begin{bmatrix} r \\ x \end{bmatrix} = \begin{bmatrix} b \\ 0 \end{bmatrix} := f, \quad (5.80)$$

此时的 KKT 矩阵是非奇异的非对称半正定矩阵.

下面将白中治等学者提出的关于非对称正定线性方程组的 HSS (Hermite/skew-Hermite Splitting) 迭代法应用于求解上述 KKT 方程组 (5.80).

选取参数 $\alpha > 0$, 将矩阵 B 分裂成

$$B = (\alpha I + H) - (\alpha I - S) = (\alpha I + S) - (\alpha I - H),$$

式中:

$$H = \frac{1}{2}(B + B^T) = \begin{bmatrix} I & O \\ O & O \end{bmatrix}, \quad S = \frac{1}{2}(B - B^T) = \begin{bmatrix} O & A \\ -A^T & O \end{bmatrix}.$$

则针对 KKT 方程组 (5.80) 的 HSS 迭代格式为

$$\begin{cases} (\alpha I + H)z^{(k+\frac{1}{2})} = (\alpha I - S)z^{(k)} + f, \\ (\alpha I + S)z^{(k+1)} = (\alpha I - H)z^{(k+\frac{1}{2})} + f. \end{cases} \quad (5.81)$$

更详细地, 即

$$\begin{cases} \begin{bmatrix} (\alpha+1)I & O \\ O & \alpha I \end{bmatrix} \begin{bmatrix} r^{(k+\frac{1}{2})} \\ x^{(k+\frac{1}{2})} \end{bmatrix} = \begin{bmatrix} \alpha I & -A \\ A^T & \alpha I \end{bmatrix} \begin{bmatrix} r^{(k)} \\ x^{(k)} \end{bmatrix} + \begin{bmatrix} b \\ 0 \end{bmatrix}, \\ \begin{bmatrix} \alpha I & A \\ -A^T & \alpha I \end{bmatrix} \begin{bmatrix} r^{(k+1)} \\ x^{(k+1)} \end{bmatrix} = \begin{bmatrix} (\alpha-1)I & O \\ O & \alpha I \end{bmatrix} \begin{bmatrix} r^{(k+\frac{1}{2})} \\ x^{(k+\frac{1}{2})} \end{bmatrix} + \begin{bmatrix} b \\ 0 \end{bmatrix}. \end{cases} \quad (5.82)$$

应用 HSS 迭代法 (5.81) 时, 求解两个子问题可采用直接法或内迭代法进行计算. 此外, 不难发现, 迭代法 (5.81) 的迭代矩阵为

$$\mathcal{L}(\alpha) = (\alpha I + S)^{-1}(\alpha I - H)(\alpha I + H)^{-1}(\alpha I - S). \quad (5.83)$$

收敛性定理如下.

定理 5.7 设 $A \in \mathbb{R}_n^{m \times n}$, $\sigma_k (k=1, 2, \dots, n)$ 为矩阵 A 的 (正) 奇异值. 则 HSS 迭代法 (5.82) 的迭代矩阵 $\mathcal{L}(\alpha)$ 有 $m-n$ 重特征值为 $\frac{\alpha-1}{\alpha+1}$, 和 $2n$ 个特征值为

$$\frac{1}{(\alpha+1)(\alpha^2+\sigma_k^2)} \left[\alpha(\alpha^2-\sigma_k^2) \pm \sqrt{(\alpha^2+\sigma_k^2)^2-4\alpha^4\sigma_k^2} \right], \quad k=1, 2, \dots, n.$$

因此, 有

$$\rho(\mathcal{L}(\alpha)) < 1, \quad \forall \alpha > 0,$$

即 HSS 迭代收敛到 KKT 方程组 (5.79) 的准确解.

证明 由于

$$\alpha I \pm H = \begin{bmatrix} (\alpha \pm 1)I & O \\ O & \alpha I \end{bmatrix}, \quad \alpha I \pm S = \begin{bmatrix} \alpha I & \pm A \\ \mp A^T & \alpha I \end{bmatrix},$$

且有

$$\alpha I + S = \begin{bmatrix} I & O \\ -\alpha^{-1}A^T & I \end{bmatrix} \begin{bmatrix} \alpha I & A \\ O & S(\alpha) \end{bmatrix},$$

这里 $S(\alpha) = \alpha I + \frac{1}{\alpha}A^T A$. 直接计算, 得

$$\mathcal{L}(\alpha) = \begin{bmatrix} \mathcal{L}_1(\alpha) & \mathcal{L}_2(\alpha) \\ \mathcal{L}_3(\alpha) & \mathcal{L}_4(\alpha) \end{bmatrix},$$

式中:

$$\begin{aligned} \mathcal{L}_1(\alpha) &= \frac{\alpha-1}{\alpha+1}I - \frac{2}{\alpha+1}AS(\alpha)^{-1}A^T, & \mathcal{L}_2(\alpha) &= -\frac{2\alpha}{\alpha+1}AS(\alpha)^{-1}, \\ \mathcal{L}_3(\alpha) &= \frac{2\alpha}{\alpha+1}S(\alpha)^{-1}A^T, & \mathcal{L}_4(\alpha) &= -\frac{\alpha-1}{\alpha+1}I + \frac{2\alpha^2}{\alpha+1}S(\alpha)^{-1}. \end{aligned}$$

设矩阵 A 的奇异值分解为

$$A = U \begin{bmatrix} \Sigma \\ O \end{bmatrix} V^T, \quad \Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_n),$$

式中: $U \in \mathbb{R}^{m \times m}$, $V \in \mathbb{R}^{n \times n}$ 为正交矩阵. 于是有

$$S(\alpha) = \alpha I + \frac{1}{\alpha}V \begin{bmatrix} \Sigma & O \end{bmatrix} U^T U \begin{bmatrix} \Sigma \\ O \end{bmatrix} V^T = V \left(\alpha I + \frac{1}{\alpha} \Sigma^2 \right) V^T,$$

及

$$\mathcal{L}_1(\alpha) = U \begin{bmatrix} \frac{\alpha-1}{\alpha+1}I - \frac{2}{\alpha+1} \left(\alpha I + \frac{1}{\alpha} \Sigma^2 \right)^{-1} \Sigma^2 & O \\ O & \frac{\alpha-1}{\alpha+1}I \end{bmatrix} U^T,$$

$$\begin{aligned}\mathcal{L}_2(\alpha) &= U \begin{bmatrix} -\frac{2\alpha}{\alpha+1} \Sigma(\alpha I + \frac{1}{\alpha} \Sigma^2)^{-1} \\ O \end{bmatrix} V^T, \\ \mathcal{L}_3(\alpha) &= V \begin{bmatrix} \frac{2\alpha}{\alpha+1} \Sigma(\alpha I + \frac{1}{\alpha} \Sigma^2)^{-1}, & O \end{bmatrix} U^T, \\ \mathcal{L}_4(\alpha) &= V \begin{bmatrix} -\frac{\alpha-1}{\alpha+1} I + \frac{2\alpha^2}{\alpha+1} (\alpha I + \frac{1}{\alpha} \Sigma^2)^{-1} \end{bmatrix} V^T.\end{aligned}$$

令

$$Q = \begin{bmatrix} U & O \\ O & V \end{bmatrix},$$

则 $\mathcal{L}(\alpha)$ 经过正交相似变换 $\tilde{\mathcal{L}}(\alpha) = Q^T \mathcal{L}(\alpha) Q$ 后, 化为

$$\tilde{\mathcal{L}}(\alpha) = \begin{bmatrix} \frac{\alpha-1}{\alpha+1} I - \frac{2}{\alpha+1} (\alpha I + \frac{1}{\alpha} \Sigma^2)^{-1} \Sigma^2 & O & -\frac{2\alpha}{\alpha+1} \Sigma(\alpha I + \frac{1}{\alpha} \Sigma^2)^{-1} \\ O & \frac{\alpha-1}{\alpha+1} I & O \\ \frac{2\alpha}{\alpha+1} \Sigma(\alpha I + \frac{1}{\alpha} \Sigma^2)^{-1} & O & -\frac{\alpha-1}{\alpha+1} I + \frac{2\alpha^2}{\alpha+1} (\alpha I + \frac{1}{\alpha} \Sigma^2)^{-1} \end{bmatrix},$$

从而可得 $\mathcal{L}(\alpha)$ 有 $m-n$ 重特征值 $\frac{\alpha-1}{\alpha+1}$ (显然其绝对值 < 1), 和矩阵

$$\begin{aligned}\mathcal{L}_k(\alpha) &= \begin{bmatrix} \frac{\alpha-1}{\alpha+1} - \frac{2}{\alpha+1} (\alpha + \frac{1}{\alpha} \sigma_k^2)^{-1} \sigma_k^2 & -\frac{2\alpha}{\alpha+1} \sigma_k (\alpha + \frac{1}{\alpha} \sigma_k^2)^{-1} \\ \frac{2\alpha}{\alpha+1} \sigma_k (\alpha + \frac{1}{\alpha} \sigma_k^2)^{-1} & -\frac{\alpha-1}{\alpha+1} + \frac{2\alpha^2}{\alpha+1} (\alpha + \frac{1}{\alpha} \sigma_k^2)^{-1} \end{bmatrix} \\ &= \frac{1}{(\alpha+1)(\alpha^2 + \sigma_k^2)} \begin{bmatrix} (\alpha-1)\alpha^2 - (\alpha+1)\sigma_k^2 & -2\alpha^2 \sigma_k \\ 2\alpha^2 \sigma_k & (\alpha+1)\alpha^2 - (\alpha-1)\sigma_k^2 \end{bmatrix}\end{aligned}$$

的特征值, $k=1, 2, \dots, n$. 而二阶矩阵 $\mathcal{L}_k(\alpha)$ 的特征方程为

$$\left(\lambda - \frac{(\alpha-1)\alpha^2 - (\alpha+1)\sigma_k^2}{(\alpha+1)(\alpha^2 + \sigma_k^2)} \right) \left(\lambda - \frac{(\alpha+1)\alpha^2 - (\alpha-1)\sigma_k^2}{(\alpha+1)(\alpha^2 + \sigma_k^2)} \right) + \frac{4\alpha^4 \sigma_k^2}{(\alpha+1)^2 (\alpha^2 + \sigma_k^2)^2} = 0,$$

化简, 得

$$\lambda^2 - \frac{2\alpha(\alpha^2 - \sigma_k^2)}{(\alpha+1)(\alpha^2 + \sigma_k^2)} \lambda + \frac{\alpha-1}{\alpha+1} = 0. \quad (5.84)$$

它的两个根为

$$\lambda_{1,2}^{(k)} = \frac{\alpha(\alpha^2 - \sigma_k^2) \pm \sqrt{(\alpha^2 + \sigma_k^2)^2 - 4\alpha^4 \sigma_k^2}}{(\alpha+1)(\alpha^2 + \sigma_k^2)}. \quad (5.85)$$

不难验证, 当 $\alpha^2 + \sigma_k^2 > 2\alpha^2 \sigma_k$ 时, 有

$$|\lambda_{1,2}^{(k)}| \leq \frac{1}{(\alpha+1)(\alpha^2 + \sigma_k^2)} \left[\alpha|\alpha^2 - \sigma_k^2| + \sqrt{(\alpha^2 + \sigma_k^2)^2 - 4\alpha^4 \sigma_k^2} \right]$$

$$= \frac{\alpha}{\alpha+1} \left[\frac{|\alpha^2 - \sigma_k^2|}{\alpha^2 + \sigma_k^2} + \sqrt{\frac{1}{\alpha^2} - \frac{4\alpha^2 \sigma_k^2}{(\alpha^2 + \sigma_k^2)^2}} \right]$$

$$< \frac{\alpha}{\alpha+1} \left(1 + \frac{1}{\alpha} \right) = 1.$$

当 $\alpha^2 + \sigma_k^2 = 2\alpha^2 \sigma_k$ 时, 有

$$|\lambda_{1,2}^{(k)}| = \frac{\alpha|\alpha^2 - \sigma_k^2|}{(\alpha+1)(\alpha^2 + \sigma_k^2)} \leq \frac{\alpha}{\alpha+1} < 1.$$

当 $\alpha^2 + \sigma_k^2 < 2\alpha^2 \sigma_k$ 时, $\lambda_{1,2}^{(k)}$ 为一对共轭复根, 此时有

$$|\lambda_{1,2}^{(k)}| = \sqrt{\lambda_1^{(k)} \lambda_2^{(k)}} = \frac{|\alpha - 1|}{\alpha + 1} < 1.$$

综合上述, $\rho(\mathcal{L}(\alpha)) < 1, \forall \alpha > 0$. 证毕. \square

此外, 还可以讨论 HSS 迭代法 (5.81) 中的最优参数 α 的选取, 详细的分析过程可参考文献 [22], 其结论如下.

定理 5.8 设 $\mathbf{A} \in \mathbb{R}^{m \times n}$ 列满秩, $\alpha > 0$ 是给定的常数. 若 $\sigma_k (k = 1, 2, \dots, n)$ 是矩阵 \mathbf{A} 的正奇异值, $\sigma_{\min} = \min_{1 \leq k \leq n} \{\sigma_k\}$, $\sigma_{\max} = \max_{1 \leq k \leq n} \{\sigma_k\}$, 则解线性方程组 (5.80) 的 HSS 迭代法的最优参数 α 为

$$\alpha^* = \arg \min_{\alpha > 0} \rho(\mathcal{L}(\alpha)) = \sqrt{\sigma_{\min} \sigma_{\max}},$$

并且

$$\rho(\mathcal{L}(\alpha^*)) = \frac{\sigma_{\max} - \sigma_{\min}}{\sigma_{\max} + \sigma_{\min}}.$$

下面考虑迭代法 (5.82) 的执行. 这一算法的计算量主要集中在求解第 2 个子方程组, 因为其系数矩阵是一个分块的反对称矩阵. 但由于其对角块的特殊性, 可以将其降阶处理. 具体迭代格式如下:

$$\begin{cases} \mathbf{r}^{(k+\frac{1}{2})} = \frac{1}{\alpha+1} (\alpha \mathbf{r}^{(k)} - \mathbf{A} \mathbf{x}^{(k)} + \mathbf{b}), \\ \mathbf{x}^{(k+\frac{1}{2})} = \mathbf{x}^{(k)} + \frac{1}{\alpha} \mathbf{A}^T \mathbf{r}^{(k)}, \\ (\alpha \mathbf{I} + \frac{1}{\alpha} \mathbf{A} \mathbf{A}^T) \mathbf{r}^{(k+1)} = (\alpha - 1) \mathbf{r}^{(k+\frac{1}{2})} - \mathbf{A} \mathbf{x}^{(k+\frac{1}{2})} + \mathbf{b}, \\ \mathbf{x}^{(k+1)} = \mathbf{x}^{(k+\frac{1}{2})} + \frac{1}{\alpha} \mathbf{A}^T \mathbf{r}^{(k+1)}. \end{cases} \quad (5.86)$$

迭代格式 (5.86) 的主要工作量集中在第 3 式, 即求解以 $\alpha \mathbf{I} + \frac{1}{\alpha} \mathbf{A} \mathbf{A}^T$ 为系数矩阵的方程组. 但由于该矩阵是对称正定的, 故可使用共轭梯度法进行非精确求解, 相应的迭代格式称为非精确 HSS 迭代法 (简记为 IHSS).

IHSS 迭代格式的 MATLAB 程序如下:

```

function [x,fval,k,time]=kkt_ihss(A,b,alpha,x,tol,max_it)
%IHSS算法求解KKT方程r=b-Ax, A'r=0
if nargin<6, max_it=1000; end
if nargin<5, tol=1.e-6; end
if nargin<4, x=zeros(size(A,2),1); end
tic; m=size(A,1); r=b-A*x;
s=A'*r; nr=norm(s); k=0;
B=alpha*eye(m)+A*A'/alpha;
while (k<=max_it)
    k=k+1;
    r=(alpha*r-A*x+b)/(alpha+1);
    x=x+s/alpha;
    c=(alpha-1)*r-A*x+b;
    r=eq_cg(B,c); %CG法求解子问题
    s=A'*r; x=x+s/alpha;
    if (norm(s)/nr<tol), break; end
end
fval=norm(r);
time=toc;

```

例 5.8 利用 HSS 迭代法的 MATLAB 程序, 计算例 5.6 中的最小二乘问题 $\min \|b - Ax\|_2$ 的极小解, 并与 CG 迭代法进行比较 (取初始点为零向量, 参数 $\alpha = 6.0$, 容许误差 $\varepsilon = 10^{-10}$).

解 编写 MATLAB 脚本文件 ex58.m, 并在命令窗口运行之, 得到数值结果如表 5.2 所示.

表 5.2 IHSS, HSS 与 CG 方法的数值结果

算法	迭代次数	CPU时间	极小解	极小值
IHSS	741	0.0003	$(-0.0303, 0.0171, 2.4508, 1.2953)^T$	0.9959
HSS	743	0.0159	$(-0.0303, 0.0171, 2.4508, 1.2953)^T$	0.9959
CG	5	0.0031	$(-0.0309, 0.0171, 2.4509, 1.2954)^T$	0.9959

习题 5

5.1 设

$$A = \begin{bmatrix} 2 & 1 \\ 1 & 0 \\ 1 & 0 \end{bmatrix}, \quad b = \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}.$$

求 LS 问题 (5.1) 的 LS 解 x_{LS} .

5.2 设

$$A = \begin{bmatrix} 1 & 3 & 1 & 1 \\ 2 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}, \quad b = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}.$$

求 LS 问题的全部解.

5.3 证明问题 (5.18) 和问题 (5.19) 的等价性.

5.4 利用等式

$$\|A(x + \alpha w) - b\|_2^2 = \|Ax - b\|_2^2 + 2\alpha w^T A^T (Ax - b) + \alpha^2 \|Aw\|_2^2.$$

证明: 若 $x \in S_{LS}$, 则 $A^T Ax = A^T b$.

5.5 给定矩阵 $A \in \mathbb{R}^{m \times n}$, 其 Moore-Penrose 广义逆矩阵记为 A^\dagger . 试证明:

- (1) 若 $\text{rank}(A) = n$, 则 $A^\dagger = (A^T A)^{-1} A^T$;
- (2) 若 $\text{rank}(A) = r \leq n$, 且 $A = BC$, 其中 $B \in \mathbb{R}^{m \times r}$, $C \in \mathbb{R}^{r \times n}$, 则 $A^\dagger = C^\dagger B^\dagger$;
- (3) $A^\dagger b$ 是超定方程组 $Ax = b$ 的极小范数最小二乘解.

5.6 设矩阵 $A \in \mathbb{R}^{m \times n}$, 且存在 $X \in \mathbb{C}^{n \times m}$ 使得对每个 $b \in \mathbb{R}^m$, $x = Xb$ 均极小化 $\|Ax - b\|_2$. 试证明: $AXA = A$ 和 $(AX)^T = AX$.

5.7 给定 $f \in \mathbb{R}^m$, $g \in \mathbb{R}^n$. 定义

$$\mathcal{X} = \{X \in \mathbb{R}^{m \times n} : X^T f = g\}.$$

证明: 当且仅当 $gf^T f = g$ 时, $\mathcal{X} \neq \emptyset$; 且当 $\mathcal{X} \neq \emptyset$ 时, 任一 $X \in \mathcal{X}$ 都可表示成

$$X = (f^T)^\dagger g^T + (I - ff^T)Z,$$

其中 $Z \in \mathbb{R}^{m \times n}$.

5.8 给定矩阵 $A \in \mathbb{C}^{m \times m}$, $B \in \mathbb{C}^{n \times n}$ 和 $C \in \mathbb{C}^{m \times n}$. 求矩阵 $X \in \mathbb{C}^{m \times n}$, 使其满足

$$\|AX - XB - C\|_F = \min,$$

并讨论 $AX - XB = C$ 相容的条件.

5.9 证明: 如果 $A \in \mathbb{R}^{m \times n}$ 的秩为 n , 则只要 $\|E\|_2 \|A^\dagger\|_2 < 1$, 就有 $A + E$ 的秩也为 n .

5.10 证明: 如果 $A_k \rightarrow A$ 和 $A_k^\dagger \rightarrow A^\dagger$ 成立, 则存在正整数 k_0 , 使得 $\forall k \geq k_0$, 秩 $\text{rank}(A_k)$ 是常数.

5.11 设 $A \in \mathbb{R}^{m \times n}$, 定义 $B(t) = (A^T A + tI)^{-1} A^T$, 其中 $t > 0$. 试证明

$$\|B(t) - A^\dagger\|_2 = \frac{t}{\sigma_r(A)[\sigma_r(A)^2 + t]}, \quad r = \text{rank}(A),$$

从而, 当 $t \rightarrow 0$ 时, 有 $B(t) \rightarrow A^\dagger$.

5.12 假定 $A^T Ax = A^T f$, $(A^T A + E)\hat{x} = A^T f$, $2\|E\|_2 \leq \sigma_n(A)^2$. 试证明: 如果 $r = f - Ax$, $\hat{r} = f - A\hat{x}$, 则有 $\hat{r} - r = A(A^T A + E)^{-1} Ex$ 和

$$\|\hat{r} - r\|_2 \leq 2\kappa_2(A) \frac{\|E\|_2}{\|A\|_2} \|x\|_2.$$

第 6 章 解线性方程组的直接法

本章考虑如下 n 阶线性方程组

$$Ax = b \quad (6.1)$$

的直接法, 其中 $A = (a_{ij})$ 称为方程组的系数矩阵, $b = (b_1, b_2, \dots, b_n)^T$ 称为方程组的右端项, 向量 $x = (x_1, x_2, \dots, x_n)^T$ 为所求的解. 若系数矩阵 A 是非奇异的, 则线性方程组 (6.1) 存在唯一解.

所谓直接法, 是指在没有舍入误差的情况下经过有限次运算可求得方程组的精确解的方法. 本章主要介绍求解线性方程组 (6.1) 最基本的直接法—Gauss 消去法和 LU 分解法以及舍入误差分析等.

6.1 Gauss 消去法

6.1.1 顺序 Gauss 消去法

Gauss 消去法的基本思想是: 首先使用初等行变换将方程组转化为一个同解的上三角形方程组 (称为消元), 再通过回代法求解该三角形方程组 (称为回代). 按行原先的位置进行消元的 Gauss 消去法称为顺序 Gauss 消去法.

例 6.1 用顺序 Gauss 消去法解线性方程组

$$\begin{cases} x_1 + x_2 + x_3 + x_4 = 6, \\ -x_1 + 2x_2 - 3x_3 + x_4 = -2, \\ 3x_1 - 3x_2 + 6x_3 - 2x_4 = 7, \\ -4x_1 + 5x_2 + 2x_3 - 3x_4 = 7. \end{cases}$$

解 (1) 消元过程:

$$\begin{bmatrix} 1 & 1 & 1 & 1 & 6 \\ -1 & 2 & -3 & 1 & -2 \\ 3 & -3 & 6 & -2 & 7 \\ -4 & 5 & 2 & -3 & 7 \end{bmatrix} \xrightarrow{\begin{matrix} r_2 + r_1 \\ r_3 - 3r_1 \\ r_4 + 4r_1 \end{matrix}} \begin{bmatrix} 1 & 1 & 1 & 1 & 6 \\ 0 & 3 & -2 & 2 & 4 \\ 0 & -6 & 3 & -5 & -11 \\ 0 & 9 & 6 & 1 & 31 \end{bmatrix} \xrightarrow{\begin{matrix} r_3 + 2r_2 \\ r_4 - 3r_2 \end{matrix}}$$

$$\begin{bmatrix} 1 & 1 & 1 & 1 & 6 \\ 0 & 3 & -2 & 2 & 4 \\ 0 & 0 & -1 & -1 & -3 \\ 0 & 0 & 12 & -5 & 19 \end{bmatrix} \xrightarrow{r_4 + 12r_3} \begin{bmatrix} 1 & 1 & 1 & 1 & 6 \\ 0 & 3 & -2 & 2 & 4 \\ 0 & 0 & -1 & -1 & -3 \\ 0 & 0 & 0 & -17 & -17 \end{bmatrix}.$$

(2) 回代过程:

$$\begin{cases} x_1 + x_2 + x_3 + x_4 = 6, \\ 3x_2 - 2x_3 + 2x_4 = 4, \\ -x_3 - x_4 = -3, \\ -17x_4 = -17. \end{cases} \Rightarrow \begin{cases} x_4 = 1, \\ x_3 = 3 - x_4 = 2, \\ x_2 = (4 + 2x_3 - 2x_4)/3 = 2, \\ x_1 = 6 - x_2 - x_3 - x_4 = 1. \end{cases}$$

故原方程组的解是: $x_1 = 1, x_2 = 2, x_3 = 2, x_4 = 1$.

对于一般线性方程组, 使用顺序 Gauss 消去法求解

$$\begin{cases} a_{11}^{(1)}x_1 + a_{12}^{(1)}x_2 + \cdots + a_{1n}^{(1)}x_n = b_1^{(1)}, \\ a_{21}^{(1)}x_1 + a_{22}^{(1)}x_2 + \cdots + a_{2n}^{(1)}x_n = b_2^{(1)}, \\ \vdots \\ a_{n1}^{(1)}x_1 + a_{n2}^{(1)}x_2 + \cdots + a_{nn}^{(1)}x_n = b_n^{(1)}. \end{cases} \quad (6.2)$$

(1) 消元过程:

$$\begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} & \cdots & a_{1n}^{(1)} & b_1^{(1)} \\ a_{21}^{(1)} & a_{22}^{(1)} & a_{23}^{(1)} & \cdots & a_{2n}^{(1)} & b_2^{(1)} \\ a_{31}^{(1)} & a_{32}^{(1)} & a_{33}^{(1)} & \cdots & a_{3n}^{(1)} & b_3^{(1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ a_{n1}^{(1)} & a_{n2}^{(1)} & a_{n3}^{(1)} & \cdots & a_{nn}^{(1)} & b_n^{(1)} \end{bmatrix} \rightarrow \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} & \cdots & a_{1n}^{(1)} & b_1^{(1)} \\ 0 & a_{22}^{(2)} & a_{23}^{(2)} & \cdots & a_{2n}^{(2)} & b_2^{(2)} \\ 0 & a_{32}^{(2)} & a_{33}^{(2)} & \cdots & a_{3n}^{(2)} & b_3^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & a_{n2}^{(2)} & a_{n3}^{(2)} & \cdots & a_{nn}^{(2)} & b_n^{(2)} \end{bmatrix}$$

$$\rightarrow \cdots \rightarrow \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} & \cdots & a_{1n}^{(1)} & b_1^{(1)} \\ 0 & a_{22}^{(2)} & a_{23}^{(2)} & \cdots & a_{2n}^{(2)} & b_2^{(2)} \\ 0 & 0 & a_{33}^{(3)} & \cdots & a_{3n}^{(3)} & b_3^{(3)} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & a_{nn}^{(n)} & b_n^{(n)} \end{bmatrix},$$

式中:

$$a_{ij}^{(2)} = a_{ij}^{(1)} - m_{i1}a_{1j}^{(1)}, \quad b_i^{(2)} = b_i^{(1)} - m_{i1}b_1^{(1)}, \quad m_{i1} = \frac{a_{i1}^{(1)}}{a_{11}^{(1)}}, \quad i, j = 2, \cdots, n.$$

一般地, 有

$$a_{ij}^{(k+1)} = a_{ij}^{(k)} - m_{ik}a_{kj}^{(k)}, \quad b_i^{(k+1)} = b_i^{(k)} - m_{ik}b_k^{(k)}, \quad m_{ik} = \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}},$$

$$i, j = k+1, \cdots, n; \quad k = 1, \cdots, n-1. \quad (6.3)$$

(2) 回代过程:

$$\begin{cases} a_{11}^{(1)}x_1 + a_{12}^{(1)}x_2 + \cdots + a_{1n}^{(1)}x_n = b_1^{(1)}, \\ a_{22}^{(2)}x_2 + \cdots + a_{2n}^{(2)}x_n = b_2^{(2)}, \\ \vdots \\ a_{nn}^{(n)}x_n = b_n^{(n)}, \end{cases} \Rightarrow \begin{cases} x_n = b_n^{(n)} / a_{nn}^{(n)}, \\ x_k = \left(b_k^{(k)} - \sum_{j=k+1}^n a_{kj}^{(k)} x_j \right) / a_{kk}^{(k)}, \quad k = n-1, \dots, 2, 1. \end{cases} \quad (6.4)$$

在此基础上, 得到顺序 Gauss 消去法的具体算法步骤.

算法 6.1 (顺序 Gauss 消去法)

步1, 输入系数矩阵 A , 右端项 b , 置 $k := 1$.

步2, 消元: 对 $k = 1, \dots, n-1$, 计算

$$\begin{aligned} m_{ik} &= a_{ik}^{(k)} / a_{kk}^{(k)}, \quad a_{ik}^{(k+1)} = 0, \\ a_{ij}^{(k+1)} &= a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)}, \quad b_i^{(k+1)} = b_i^{(k)} - m_{ik} b_k^{(k)}. \\ (i &= k+1, \dots, n; \quad j = k+1, \dots, n.) \end{aligned}$$

步3, 回代:

$$\begin{aligned} x_n &= b_n^{(n)} / a_{nn}^{(n)}, \\ x_k &= \left(b_k^{(k)} - \sum_{j=k+1}^n a_{kj}^{(k)} x_j \right) / a_{kk}^{(k)}, \quad k = n-1, \dots, 1. \end{aligned}$$

可以统计顺序 Gauss 消去法的计算量. 由于加减法的计算量可忽略不计, 此处只统计乘除法次数.

消元过程: 第 k ($k = 1, \dots, n-1$) 步消元, 有

$$(n-k)(n-k+1) + (n-k) = (n-k)(n-k+2)$$

次乘除法, 共

$$\begin{aligned} N_1 &= \sum_{k=1}^{n-1} (n-k)(n-k+2) = \sum_{i=1}^{n-1} (i^2 + 2i) \\ &= \frac{n(n-1)(2n-1)}{6} + n(n-1) = \frac{n(n-1)(2n+5)}{6} \end{aligned}$$

次乘除法.

回代过程: 计算 x_k ($k = n, \dots, 2, 1$) 时, 有 $n-k+1$ 次乘除法, 共

$$N_2 = \sum_{k=1}^n (n-k+1) = \sum_{i=1}^n i = \frac{n(n+1)}{2}$$

次乘法.

消元和回代过程共计

$$N_1 + N_2 = \frac{n(n-1)(2n+5)}{6} + \frac{n(n+1)}{2} = \frac{n^3}{3} + n^2 - \frac{n}{3}$$

次乘法.

可见消元过程的计算量为 $O(n^3)$, 而回代过程的计算量为 $O(n^2)$, 因此顺序 Gauss 消去法的计算量主要在消元过程部分.

算法 6.1 要求对所有的 $k = 1, \dots, n$, $a_{kk}^{(k)} \neq 0$, 此时称顺序 Gauss 消去法是可行的. 一般将 $a_{kk}^{(k)}$ 称为第 k 步消元的主元. 下面的定理给出了顺序 Gauss 消去法可行的一个充要条件.

定理 6.1 证明: 顺序 Gauss 消去法可行的充分必要条件是系数矩阵 A 的所有顺序主子式 $D_k \neq 0$, $k = 1, 2, \dots, n$.

证明 必要性. 若顺序 Gauss 消去法是可行的, 即 $a_{kk}^{(k)} \neq 0$, 则可进行消去法的 $k-1$ 步 ($k \leq n$). 由于 $A^{(k)}$ 是由 A 逐行实行初等变换 (某数乘以某一行加到另一行) 得到的, 这些运算不改变相应顺序主子式的值, 故有

$$D_k = \begin{vmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \cdots & a_{1k}^{(1)} \\ & a_{22}^{(2)} & \cdots & a_{2k}^{(2)} \\ & & \ddots & \vdots \\ & & & a_{kk}^{(k)} \end{vmatrix} = a_{11}^{(1)} a_{22}^{(2)} \cdots a_{kk}^{(k)} \neq 0, \quad k = 1, 2, \dots, n.$$

充分性. 用归纳法证明. 当 $k = 1$ 时显然成立. 设命题对 $k-1$ 成立. 现设 $D_1 \neq 0, \dots, D_{k-1} \neq 0, D_k \neq 0$. 由归纳法假设有 $a_{11}^{(1)} \neq 0, \dots, a_{k-1,k-1}^{(k-1)} \neq 0$. 因此, 消去法可以进行第 $k-1$ 步, A 约化为

$$A^{(k)} = \begin{bmatrix} A_{11}^{(k-1)} & A_{12}^{(k-1)} \\ & A_{22}^{(k)} \end{bmatrix},$$

式中: $A_{11}^{(k-1)}$ 是对角元为 $a_{11}^{(1)}, \dots, a_{k-1,k-1}^{(k-1)}$ 的上三角矩阵. 因 $A^{(k)}$ 是通过行初等变换由 A 逐步得到的, 故 A 的 k 阶顺序主子式与 $A^{(k)}$ 的 k 阶顺序主子式相等, 即

$$D_k = \det \begin{bmatrix} A_{11}^{(k-1)} & \alpha_{12}^{(k-1)} \\ & a_{kk}^{(k)} \end{bmatrix} = a_{11}^{(1)} \cdots a_{k-1,k-1}^{(k-1)} a_{kk}^{(k)},$$

式中: $\alpha_{12}^{(k-1)}$ 为 $A_{12}^{(k-1)}$ 的第 1 列. 故由 $D_k \neq 0$ 及归纳法假设可推出 $a_{kk}^{(k)} \neq 0$. 证毕. \square

根据算法 6.1, 可编制 MATLAB 程序如下:

```
%顺序Gauss消去法程序-msgauss.m
function x=msgauss(A,b,flag)
```



```

%输入:A为系数矩阵,b为右端项,若flag=0(默认),
%    则不显示中间过程,否则显示中间过程
%输出:x为解向量
if nargin<3, flag=0; end
n=length(b);
for k=1:(n-1) %消元过程
    m=A(k+1:n,k)/A(k,k);
    A(k+1:n,k+1:n)=A(k+1:n,k+1:n)-m*A(k,k+1:n);
    b(k+1:n)=b(k+1:n)-m*b(k);
    A(k+1:n,k)=zeros(n-k,1);
    if flag~=0, Ab=[A,b], end
end
x=zeros(n,1); %回代过程
x(n)=b(n)/A(n,n);
for k=n-1:-1:1
    x(k)=(b(k)-A(k,k+1:n)*x(k+1:n))/A(k,k);
end

```

例 6.2 利用程序 msgauss.m 计算下列方程组的解

$$\begin{cases} x_1 + x_2 + x_3 + x_4 = 6, \\ -x_1 + 2x_2 - 3x_3 + x_4 = -2, \\ 3x_1 - 3x_2 + 6x_3 - 2x_4 = 7, \\ -4x_1 + 5x_2 + 2x_3 - 3x_4 = 7. \end{cases}$$

解 在 MATLAB 命令窗口输入:

```

>> A=[1 1 1 1;-1 2 -3 1;3 -3 6 -2;-4 5 2 -3];
>> b=[6 -2 7 7]';
>> x=msgauss(A,b)

```

可得计算结果 (略).

6.1.2 列主元 Gauss 消去法

一般来说, 顺序 Gauss 消去法的计算过程是不可靠的, 一旦出现 $a_{kk}^{(k)} = 0$, 计算就无法进行下去. 另外, 即使对所有 $k = 1, 2, \dots, n$, $a_{kk}^{(k)} \neq 0$, 也不能保证计算过程是数值稳定的.

例 6.3 设有线性方程组

$$\begin{cases} 0.0001x_1 + 1.0x_2 = 1.0, \\ 1.0x_1 + 1.0x_2 = 2.0. \end{cases}$$

其精确解为

$$x_1 = \frac{10000}{9999} \approx 1.00010, \quad x_2 = \frac{9998}{9999} \approx 0.99990.$$

现在假定用尾数为 4 位十进制字长的浮点数来求解.

解 (1) 消元过程: 根据 4 位浮点数运算规则 $1.0 - 10000.0 = (0.00001 - 0.1)10^5 = (0.0000 - 0.1)10^5 = -10000.0$ (舍入). 同理, $2.0 - 10000.0 = -10000.0$.

$$\begin{aligned} & \left[\begin{array}{ccc} 0.0001 & 1.0 & 1.0 \\ 1.0 & 1.0 & 2.0 \end{array} \right] \xrightarrow{r_2 - 10^4 r_1} \left[\begin{array}{ccc} 0.0001 & 1.0 & 1.0 \\ 0 & 1.0 - 10000.0 & 2.0 - 10000.0 \end{array} \right] \\ & \xrightarrow{\text{舍入}} \left[\begin{array}{ccc} 0.0001 & 1.0 & 1.0 \\ 0 & -10000.0 & -10000.0 \end{array} \right]. \end{aligned}$$

(2) 回代过程:

$$\begin{cases} 0.0001x_1 + 1.0x_2 = 1.0, \\ -10000.0x_2 = -10000.0. \end{cases} \Rightarrow \begin{cases} x_2 = 1.0, \\ x_1 = 0.0. \end{cases}$$

代入原方程组验算, 发现结果严重失真.

分析结果失真的原因发现, 由于第 1 列的主元素 0.0001 绝对值过于小, 当它在消元过程中作分母时把中间过程数据放大 10000 倍, 使中间结果“吃”掉了原始数据, 从而造成数值不稳定.

针对以上问题, 考虑选用绝对值大的数作为主元素.

(1) 消元过程:

$$\begin{aligned} & \left[\begin{array}{ccc} 0.0001 & 1.0 & 1.0 \\ 1.0 & 1.0 & 2.0 \end{array} \right] \xrightarrow{r_1 \leftrightarrow r_2} \left[\begin{array}{ccc} 1.0 & 1.0 & 2.0 \\ 0.0001 & 1.0 & 1.0 \end{array} \right] \\ & \xrightarrow{r_2 - 0.0001r_1} \left[\begin{array}{ccc} 1.0 & 1.0 & 2.0 \\ 0 & 1.0 - 0.0001 & 1.0 - 0.0002 \end{array} \right] \\ & \xrightarrow{\text{舍入}} \left[\begin{array}{ccc} 1.0 & 1.0 & 2.0 \\ 0 & 1.0 & 1.0 \end{array} \right]. \end{aligned}$$

这里, 舍入过程 $1.0 - 0.0001 = (0.1 - 0.00001)10^1$ (舍入), 同理 $1.0 - 0.0002 = 1.0$.

(2) 回代过程:

$$\begin{cases} 1.0x_1 + 1.0x_2 = 2.0, \\ 1.0x_2 = 1.0. \end{cases} \Rightarrow \begin{cases} x_2 = 1.0, \\ x_1 = 1.0. \end{cases}$$

代入原方程组验算, 发现结果基本合理.

例 6.3 说明了选主元素的重要性. 下面阐述列主元 Gauss 消去法的基本思想. 记 $A^{(1)} = A$, 在消元过程的第 1 步, 取第 1 列中绝对值最大的元素 $a_{r_1 1}^{(1)}$, 即

$$a_{r_1 1}^{(1)} = \max_{1 \leq i \leq n} |a_{i1}^{(1)}|$$

作为主元素. 若 $r_1 > 1$, 交换第 r_1 行和第 1 行.

一般地, 在消元过程的第 k 步, 取

$$a_{r_k k}^{(k)} = \max_{k \leq i \leq n} |a_{ik}^{(k)}| \quad (6.5)$$

作为主元素. 若 $r_k > k$, 交换第 r_k 行和第 k 行.

列主元 Gauss 消去法的算法具体步骤如下.

算法 6.2 (列主元 Gauss 消去法)

步1, 输入系数矩阵 A , 右端项 b , 置 $k := 1$.

步2, 对 $k = 1, \dots, n-1$ 进行如下操作:

① 选列主元, 确定 r_k , 使

$$a_{r_k k}^{(k)} = \max_{k \leq i \leq n} |a_{ik}^{(k)}|,$$

若 $a_{r_k k}^{(k)} = 0$, 则停止计算, 否则, 进行下一步.

② 若 $r_k > k$, 交换 $(A^{(k)}, b^{(k)})$ 的第 k, r_k 两行.

③ 消元: 对 $i, j = k+1, \dots, n$, 计算

$$\begin{aligned} m_{ik} &= a_{ik}^{(k)} / a_{kk}^{(k)}, \quad a_{ik}^{(k+1)} = 0, \\ a_{ij}^{(k+1)} &= a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)}, \quad b_i^{(k+1)} = b_i^{(k)} - m_{ik} b_k^{(k)}. \end{aligned}$$

步3, 回代:

$$\begin{aligned} x_n &= b_n^{(n)} / a_{nn}^{(n)}, \\ x_k &= \left(b_k^{(k)} - \sum_{j=k+1}^n a_{kj}^{(k)} x_j \right) / a_{kk}^{(k)}, \quad k = n-1, \dots, 1. \end{aligned}$$

根据算法 6.2, 编制 MATLAB 程序如下:

```
%列主元Gauss消去法程序-mgauss2.m
function x=mgauss2(A,b,flag)
%输入:A为系数矩阵,b为右端项,若flag=0(默认),
%      则不显示中间过程,否则显示中间过程
%输出:x为解向量
if nargin<3, flag=0; end
n=length(b);
for k=1:(n-1) %选主元
    [ap,p]=max(abs(A(k:n,k)));p=p+k-1;
    if p>k
        A([k p],:)=A([p k],:);b([k p],:)=b([p k],:);
```

```

end
m=A(k+1:n,k)/A(k,k); %消元
A(k+1:n,k+1:n)=A(k+1:n,k+1:n)-m*A(k,k+1:n);
b(k+1:n)=b(k+1:n)-m*b(k);A(k+1:n,k)=zeros(n-k,1);
if flag~=0, Ab=[A,b], end
end
x=zeros(n,1); x(n)=b(n)/A(n,n); %回代
for k=n-1:-1:1
    x(k)=(b(k)-A(k,k+1:n)*x(k+1:n))/A(k,k);
end

```

例 6.4 利用程序 mgauss2.m 计算下列线性方程组的解

$$\begin{bmatrix} 2 & -1 & 4 & -3 & 1 \\ -1 & 1 & 2 & 1 & 3 \\ 4 & 2 & 3 & 3 & -1 \\ -3 & 1 & 3 & 2 & 4 \\ 1 & 3 & -1 & 4 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} = \begin{bmatrix} 11 \\ 14 \\ 4 \\ 16 \\ 18 \end{bmatrix}.$$

解 在 MATLAB 命令窗口输入:

```

>> A=[2 -1 4 -3 1;-1 1 2 1 3;4 2 3 3 -1;-3 1 3 2 4;1 3 -1 4 4];
>> b=[11 14 4 16 18]';
>> x=mpgauss(A,b); x=x'

```

得计算结果:

```

x =
    1.0000    2.0000    1.0000   -1.0000    4.0000

```

6.2 LU 分解法

把一个 n 阶矩阵分解成两个三角形矩阵的乘积称为矩阵的三角分解. 本节介绍矩阵的 LU 分解 $A = LU$, 其中 L 是单位下三角矩阵, U 是上三角矩阵. 这种形式的分解对于求解方程组 (6.1) 是十分有用的. 事实上, 若 $A = LU$ 是一个 LU 分解, 此时线性方程组

$$Ax = b \Rightarrow LUx = b \Rightarrow \begin{cases} Ly = b \\ Ux = y \end{cases} \quad (6.6)$$

转化为 $Ly = b$ 及 $Ux = y$ 两个三角形方程组. 由于三角形方程组很容易通过向前消去法或回代方法求解, 且只有 $O(n^2)$ 的计算量, 故研究矩阵的 LU 分解十分有意义.

6.2.1 顺序 LU 分解法

首先讨论矩阵的 LU 分解. 设 $A = LU$, 其中 L 为一个单位下三角矩阵, U 为一个上三角矩阵, 即

$$L = \begin{bmatrix} 1 & & & \\ l_{21} & 1 & & \\ \vdots & \vdots & \ddots & \\ l_{n1} & l_{n2} & \cdots & 1 \end{bmatrix}, \quad U = \begin{bmatrix} u_{11} & u_{12} & \cdots & u_{1n} \\ & u_{22} & \cdots & u_{2n} \\ & & \ddots & \vdots \\ & & & u_{nn} \end{bmatrix}. \quad (6.7)$$

下面推导三角形矩阵 L 和 U 的元素的计算公式. 由等式 $A = LU$, 得

$$a_{ij} = \begin{bmatrix} l_{i1}, & \cdots, & l_{i,i-1}, & 1, & 0, & \cdots, & 0 \end{bmatrix} \begin{bmatrix} u_{1j} \\ \vdots \\ u_{j-1,j} \\ u_{jj} \\ 0 \\ \vdots \\ 0 \end{bmatrix}. \quad (6.8)$$

当 $j \geq i$ 时, 有

$$a_{ij} = l_{i1}u_{1j} + \cdots + l_{i,i-1}u_{i-1,j} + u_{ij},$$

于是

$$u_{ij} = a_{ij} - \sum_{r=1}^{i-1} l_{ir}u_{rj};$$

当 $j < i$ 时, 有

$$a_{ij} = l_{i1}u_{1j} + \cdots + l_{i,j-1}u_{j-1,j} + l_{ij}u_{jj},$$

于是

$$l_{ij} = \left(a_{ij} - \sum_{r=1}^{j-1} l_{ir}u_{rj} \right) / u_{jj}.$$

即

$$u_{1j} = a_{1j}, \quad j = 1, \cdots, n; \quad l_{i1} = a_{i1}/u_{11}, \quad i = 2, \cdots, n; \quad (6.9)$$

$$u_{ij} = a_{ij} - \sum_{r=1}^{i-1} l_{ir}u_{rj}, \quad i = 2, \cdots, n; \quad j = i, \cdots, n; \quad (6.10)$$

$$l_{ij} = \left(a_{ij} - \sum_{r=1}^{j-1} l_{ir}u_{rj} \right) / u_{jj}, \quad i = 3, \cdots, n; \quad j = 2, \cdots, i-1. \quad (6.11)$$

为了便于编程计算, 将式 (6.9) 第 1 式中的下标 j 换成 i , 式 (6.10) 中的下标 i 换成 k , 下标 j 换成 i , 式 (6.11) 中的下标 j 换成 k , 则有

$$u_{1i} = a_{1i}, \quad i = 1, \cdots, n; \quad l_{i1} = a_{i1}/u_{11}, \quad i = 2, \cdots, n; \quad (6.12)$$

$$u_{ki} = a_{ki} - \sum_{r=1}^{k-1} l_{kr} u_{ri}, \quad k = 2, \cdots, n; \quad i = k, \cdots, n; \quad (6.13)$$

$$l_{ik} = \left(a_{ik} - \sum_{r=1}^{k-1} l_{ir} u_{rk} \right) / u_{kk}, \quad k = 2, \cdots, n-1; \quad i = k+1, \cdots, n. \quad (6.14)$$

上述这种按矩阵 \mathbf{A} 元素的自然顺序进行分解的方法称为顺序 LU 分解法. 下面是用顺序 LU 分解求解线性方程组的算法步骤.

算法 6.3 (顺序 LU 分解法)

步1, 输入系数矩阵 \mathbf{A} , 右端项 \mathbf{b} .

步2, LU 分解:

$$u_{1i} = a_{1i}, \quad i = 1, \cdots, n;$$

$$l_{i1} = a_{i1}/u_{11}, \quad i = 2, \cdots, n;$$

对 $k = 2, \cdots, n$, 计算

$$u_{ki} = a_{ki} - \sum_{r=1}^{k-1} l_{kr} u_{ri}, \quad i = k, \cdots, n;$$

$$l_{ik} = (a_{ik} - \sum_{r=1}^{k-1} l_{ir} u_{rk}) / u_{kk}, \quad i = k+1, \cdots, n.$$

步3, 用向前消去法解下三角方程组 $\mathbf{L}\mathbf{y} = \mathbf{b}$:

$$y_1 = b_1; \quad y_k = b_k - \sum_{i=1}^{k-1} l_{ki} y_i, \quad k = 2, \cdots, n.$$

步4, 用回代法解上三角方程组 $\mathbf{U}\mathbf{x} = \mathbf{y}$:

$$x_n = y_n / u_{nn}; \quad x_k = \left(y_k - \sum_{i=k+1}^n u_{ki} x_i \right) / u_{kk}, \quad k = n-1, \cdots, 1.$$

注 6.1 可以看出, 利用 LU 分解, 分开了系数矩阵的计算和对右端项的计算. 正是这一特点, 使得 LU 分解法特别适用于求解系数矩阵相同而右端项不同的一系列方程组, 而控制论等领域中刚好存在这样的实际问题.

下面考虑顺序 LU 分解的程序实现. 注意到 LU 分解后, 原系数矩阵 \mathbf{A} 的数据不再需要保留, 因此, 为了节省存储空间, 可在 MATLAB 程序中将分解后的单位下三角矩阵 \mathbf{L} 和上三角矩阵 \mathbf{U} 分别存放在系数矩阵 \mathbf{A} 的严格下三角和上三角部分 (单位下三角矩阵的对角线元素 1 不需存储), 而不再为其开辟额外的存储单元.

算法 6.3 的 MATLAB 程序如下:

```
%顺序LU分解法程序-mslu.m
function [x,A]=mslu(A,b)
```

```

%输入:A为系数矩阵,b为右端向量
%输出:x为解向量,L和U分别存放在A的严格下三角和上三角部分
n=length(b);
for k=1:n %顺序LU分解
    A(k:n,k)=A(k:n,k)-A(k:n,1:k-1)*A(1:k-1,k);
    A(k+1:n,k)=A(k+1:n,k)/A(k,k); %乘子向量
    A(k,k+1:n)=A(k,k+1:n)-A(k,1:k-1)*A(1:k-1,k+1:n);
end
y=zeros(n,1);
for k=1:n, %解下三角矩阵Ly=b
    y(k)=b(k)-A(k,1:k-1)*y(1:k-1);
end
x=zeros(n,1);
for k=n:-1:1, %解上三角方程组Ux=y
    x(k)=(y(k)-A(k,k+1:n)*x(k+1:n))/A(k,k);
end

```

例 6.5 利用程序 mslu.m 计算例 6.4 中的线性方程组的解.

解 在 MATLAB 命令窗口输入:

```

>> A=[2 -1 4 -3 1;-1 1 2 1 3;4 2 3 3 -1;-3 1 3 2 4;1 3 -1 4 4];
>> b=[11 14 4 16 18]';
>> [x,A]=mslu(A,b)

```

可得计算结果 (略).

6.2.2 列主元 LU 分解法

顺序 LU 分解法在本质上与顺序 Gauss 消去法是一致的. 由算法 6.1 可知, 顺序 Gauss 消去法的第 1 步消元相当于用矩阵

$$M_1 = \begin{bmatrix} 1 & & & & \\ -m_{21} & 1 & & & \\ -m_{31} & & 1 & & \\ \vdots & & & \ddots & \\ -m_{n1} & & & & 1 \end{bmatrix}$$

左乘 $[A^{(1)}, b^{(1)}]$, 这里 m_{i1} 由式 (6.3) 所定义, 即

$$[A^{(2)}, b^{(2)}] = M_1[A^{(1)}, b^{(1)}].$$

第 2 步消元相当于用矩阵

$$M_2 = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & -m_{32} & 1 & \\ & \vdots & & \ddots \\ & -m_{n2} & & & 1 \end{bmatrix}$$

左乘 $[A^{(2)}, b^{(2)}]$, 即

$$[A^{(3)}, b^{(3)}] = M_2[A^{(2)}, b^{(2)}] = M_2 M_1[A^{(1)}, b^{(1)}].$$

一般地, 第 k 步相当于用矩阵

$$M_k = \begin{bmatrix} 1 & & & & \\ & \ddots & & & \\ & & 1 & & \\ & & -m_{k+1,k} & 1 & \\ & & \vdots & & \ddots \\ & & -m_{nk} & & & 1 \end{bmatrix}$$

左乘 $[A^{(k)}, b^{(k)}]$, 即

$$\begin{aligned} [A^{(k+1)}, b^{(k+1)}] &= M_k[A^{(k)}, b^{(k)}] = \cdots \\ &= M_k \cdots M_2 M_1[A^{(1)}, b^{(1)}], \quad k = 1, \cdots, n-1. \end{aligned}$$

由顺序 Gauss 消去法可知, 经过 $n-1$ 步消元后, 系数矩阵 A 被化成了上三角矩阵, 即 $A^{(n)} = U$, 从而

$$U = A^{(n)} = M_{n-1}A^{(n-1)} = M_{n-1} \cdots M_2 M_1 A,$$

于是

$$A = M_1^{-1} M_2^{-1} \cdots M_{n-1}^{-1} U = LU,$$

式中:

$$L = M_1^{-1} M_2^{-1} \cdots M_{n-1}^{-1} = \begin{bmatrix} 1 & & & & \\ m_{21} & 1 & & & \\ \vdots & & \ddots & & \\ m_{n-1,1} & m_{n-1,2} & \cdots & 1 & \\ m_{n1} & m_{n2} & \cdots & m_{n,n-1} & 1 \end{bmatrix}$$

为单位下三角矩阵. 由此可见, 顺序 Gauss 消去法实际上就是将方程组的系数矩阵分解成单位下三角矩阵与上三角矩阵的乘积. 对比算法 6.1 和算法 6.3, 不难看出, 顺序 Gauss 消去法的消元过程相当于 LU 分解过程和 $Ly = b$ 的求解, 而回代过程则相当于解线性方程组 $Ux = y$.

定理 6.2 若 n 阶方阵 A 的所有顺序主子式都不等于零, 则 A 存在唯一的 LU 分解 $A = LU$.

证明 因顺序 LU 分解本质上等同于顺序 Gauss 消去法, 故存在性由定理 6.1 立即可得. 下面证明唯一性. 事实上, 若 A 存在两种不同的三角分解:

$$A = LU = L_1 U_1,$$

式中: L 和 L_1 均为单位下三角矩阵; U 和 U_1 均为上三角矩阵. 因 A 是非奇异的, 故 U 和 U_1 也是非奇异的. 于是由上式, 得

$$L_1^{-1} L = U_1 U^{-1}.$$

注意到上式的左边是单位下三角矩阵, 而右边则是上三角矩阵, 故必有

$$L_1^{-1} L = U_1 U^{-1} = I \text{ (单位阵),}$$

即 $L_1 = L, U_1 = U$. 证毕. □

根据上面的分析, 既然顺序 LU 分解法本质上等同于顺序 Gauss 消去法, 因此, 为了提高计算的数值稳定性, 有必要考虑列主元 LU 分解技术. 这只需要在一般 LU 分解的第 k 步避免绝对值较小的 u_{kk} 作除数即可. 假设第 $k-1$ 步已经完成, 在进行第 k 步分解之前进行选主元的操作. 可以引入

$$s_i = a_{ik} - \sum_{r=1}^{k-1} l_{ir} u_{rk}, \quad i = k, k+1, \dots, n,$$

且令

$$|s_{i_k}| = \max_{k \leq i \leq n} |s_i|.$$

然后用 s_{i_k} 作为 u_{kk} 并交换增广矩阵 $[A, b]$ 的第 k 行和第 i_k 行, 于是有 $|l_{ik}| \leq 1$ ($i = k+1, \dots, n$), 再进行第 k 步分解. 算法如下:

算法 6.4 (列主元 LU 分解法)

步1, 输入系数矩阵 A , 右端项 b .

步2, 列主元 LU 分解:

对 $k = 1, \dots, n$,

① 计算 $s_i = a_{ik} - \sum_{r=1}^{k-1} l_{ir} u_{rk} \Rightarrow a_{ik}, \quad i = k, k+1, \dots, n.$

② 选主元 $|s_{i_k}| = \max_{k \leq i \leq n} |s_i|$, 并记录 $i_k, s_{i_k} \Rightarrow u_{kk}.$

③ 交换 $[A, b]$ 的第 k 行和第 i_k 行元素.

④ 计算 L 的第 k 列元素: $l_{ik} = s_i / u_{kk} = a_{ik} / a_{kk} \Rightarrow a_{ik}, \quad i = k+1, \dots, n.$

⑤ 计算 U 的第 k 行元素: $u_{kj} = a_{kj} - \sum_{r=1}^{k-1} l_{kr} u_{rj} \Rightarrow a_{kj}, \quad j = k+1, \dots, n.$

步3, 用向前消去法解下三角方程组 $Ly = b$:

$$y_1 = b_1; y_k = b_k - \sum_{j=1}^{k-1} l_{kj} y_j, \quad k = 2, \dots, n.$$

步4, 用回代法解上三角方程组 $Ux = y$:

$$x_n = y_n / u_{nn}; x_k = (y_k - \sum_{j=k+1}^n u_{kj} x_j) / u_{kk}, \quad k = n-1, \dots, 1.$$

下面给出列主元 LU 分解法的 MATLAB 程序:

```
%列主元LU分解法程序-mplu.m
function [x,A,P]=mplu(A,b)
%列主元LU分解PA=LU
%输入:A为系数矩阵,b为右端向量
%输出:x返回解向量,L和U分别存放在A的严格下三角和
%    上三角部分,P返回选主元时记录行交换的置换阵
n=length(b);
P=eye(n); %P记录选择主元时候所进行的行变换
for k=1:n %列主元LU分解
    A(k:n,k)=A(k:n,k)-A(k:n,1:k-1)*A(1:k-1,k);
    [s,m]=max(abs(A(k:n,k))); %选列主元
    m=m+k-1;
    if m~=k
        A([k m],:)=A([m k],:);P([k m],:)=P([m k],:);
    end
    A(k+1:n,k)=A(k+1:n,k)/A(k,k);
    A(k,k+1:n)=A(k,k+1:n)-A(k,1:k-1)*A(1:k-1,k+1:n);
end
b=P*b; y=zeros(n,1);
for k=1:n, %解下三角矩阵Ly=b
    y(k)=b(k)-A(k,1:k-1)*y(1:k-1);
end
x=zeros(n,1);
for k=n:-1:1, %解上三角方程组Ux=y
    x(k)=(y(k)-A(k,k+1:n)*x(k+1:n))/A(k,k);
end
```

例 6.6 利用程序 mplu.m 计算下列线性方程组的解

$$\begin{bmatrix} 2 & -1 & 4 & -3 & 1 \\ -1 & 1 & 2 & 1 & 3 \\ 4 & 2 & 3 & 3 & -1 \\ -3 & 1 & 3 & 2 & 4 \\ 1 & 3 & -1 & 4 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} = \begin{bmatrix} 11 \\ 14 \\ 4 \\ 16 \\ 18 \end{bmatrix}.$$

解 在 MATLAB 命令窗口输入:

```
>> A=[2 -1 4 -3 1;-1 1 2 1 3;4 2 3 3 -1;-3 1 3 2 4;1 3 -1 4 4];
>> b=[11 14 4 16 18]';
>> [x,A,P]=mplu(A,b)
```

可得计算结果 (略).

6.2.3 不完全 LU 分解

当 A 的非零元较少且按一定规则分布时 (称为非零元的稀疏模式), 其 LU 分解 (又称完全 LU 分解) 产生的单位下三角矩阵 L 和上三角矩阵 U 一般不能保持和 A 相同的稀疏模式. 本节讨论矩阵 A 的不完全 LU 分解, 即在保持稀疏模式的前提下进行 LU 分解, 这样分解后 L 和 U 的乘积是 A 的一个近似: $A \approx LU$. 这种分解得到的 L 和 U 的逆矩阵往往用来对方程组 $Ax = b$ 作预处理:

$$\tilde{A}\tilde{x} = \tilde{b}, \quad \tilde{A} = L^{-1}AU^{-1}, \quad \tilde{x} = Ux, \quad \tilde{b} = L^{-1}b.$$

使得预处理后方程组的系数矩阵有更好的条件数.

定义 6.1 设 $A = (a_{ij}) \in \mathbb{R}^{n \times n}$, 非对角元的位置和除对角线之外的其他非零元的位置分别记为下面的两个集合

$$ND = \{(i, j) : 1 \leq i, j \leq n \text{ 且 } i \neq j\}, \quad NZ = \{(i, j) \in ND : a_{ij} \neq 0\}.$$

设指标集 \mathcal{F} 满足 $NZ \subseteq \mathcal{F} \subseteq ND$, $0 \leq \omega \leq 1$. 若矩阵 A 有分解

$$A = LU + R,$$

式中: $L = (l_{ij}) \in \mathbb{R}^{n \times n}$ 为单位下三角阵且满足

$$l_{ij} = 0, \quad i > j, \quad (i, j) \notin \mathcal{F}; \quad (6.15)$$

$U = (u_{ij}) \in \mathbb{R}^{n \times n}$ 为上三角阵且满足

$$u_{ij} = 0, \quad i < j, \quad (i, j) \notin \mathcal{F}; \quad (6.16)$$

矩阵 $R = (r_{ij}) \in \mathbb{R}^{n \times n}$ 满足

$$r_{ij} = 0, \quad (i, j) \in \mathcal{F}, \quad (6.17)$$

$$r_{ii} = -\omega \sum_{j=1, j \neq i}^n r_{ij}, \quad i = 1, 2, \dots, n. \quad (6.18)$$

则称 A 有关于 \mathcal{F} 和 ω 的松弛不完全 LU 分解, ω 称为松弛参数.

若 $\mathcal{F} = ND$, 由式 (6.17) 表明 R 的非对角元为 0, 式 (6.18) 表明 R 的对角元也为 0, 所以 $R = O$. 同时式 (6.15) 和式 (6.16) 也不需要满足, 此时的松弛不完全 LU 分解

即为完全 LU 分解. 若 $\mathcal{F} = \text{NZ}$, 式 (6.17) 和式 (6.18) 表明 L 的下三角部分以及 U 的上三角部分和 A 的零元分布是一样的. 在分解中可以把集合

$$\text{ND} \setminus \mathcal{F} = \{(i, j) : 1 \leq i \neq j \leq n, (i, j) \notin \mathcal{F}\}$$

看成矩阵 A 的零模式, 而 L 和 U 保持了这种零模式. 上述松弛不完全 LU 分解中, 若参数 $\omega = 0$, 相应的分解称为不完全 LU 分解 (ILU). 若参数 $\omega = 1$, 相应的分解称为修正的不完全 LU 分解 (MILU).

下面给出矩阵 A 关于 \mathcal{F} 和 ω 的松弛不完全 LU 分解, 该分解可以通过 $n-1$ 步 Gauss 消去法和 $n-1$ 次修正来实现, 称为修正的 Gauss 消去法. 具体步骤如下: 记 $A^{(1)} = (a_{ij}^{(1)}) = A$, 假定已经执行了 $k-1$ 步修正的 Gauss 消去法并产生了 $A^{(k)} = (a_{ij}^{(k)})$, 则第 k 步修正的 Gauss 消去法如下.

(1) Gauss 消去法. 把 $A^{(k)}$ 的第 k 列中第 $k+1$ 个元素到第 n 个元素变为 0, 令

$$\tilde{A}^{(k)} = L_k A^{(k)} = A^{(k)} - l_k a_k,$$

式中: $L_k = I - l_k e_k^T$, 且

$$l_k = (0, \dots, 0, l_{k+1,k}, \dots, l_{nk})^T, \quad l_{ik} = a_{ik}^{(k-1)} / a_{kk}^{(k-1)}, \quad i = k+1, \dots, n,$$

$$a_k = e_k^T A^{(k)} = (0, \dots, 0, a_{k,k}^{(k)}, \dots, a_{kn}^{(k)})^T.$$

(2) 修正 $\tilde{A}^{(k)} = (\tilde{a}_{ij}^{(k)})$ 右下角的 $n-k$ 阶子矩阵. 将 $\tilde{A}^{(k-1)}$ 修正为

$$A^{(k+1)} = \tilde{A}^{(k)} + R_k,$$

式中: $R_k = (r_{ij}^{(k)})$ 为

$$r_{ij}^{(k)} = \begin{cases} \tilde{a}_{ij}^{(k-1)}, & k+1 \leq i, j \leq n, i \neq j, (i, j) \notin \mathcal{F}; \\ -\omega \sum_{\substack{l=k+1 \\ l \neq i, (i,l) \notin \mathcal{F}}} \tilde{a}_{il}^{(k-1)}, & k+1 \leq i, j \leq n, i = j; \\ 0, & \text{其他.} \end{cases}$$

第 $n-1$ 步修正的 Gauss 消去法后, 可得单位下三角矩阵

$$L = (L_{n-1} L_{n-2} \cdots L_1)^{-1}$$

和上三角阵 $U = A^{(n)}$. 令

$$R = \sum_{k=1}^{n-1} R_k, \quad R_k = (r_{ij}^{(k)}),$$

则可得如下定理 (见文献 [13]).

定理 6.3 设 $a_{kk}^{(k)} \neq 0, k = 1, 2, \dots, n$, 则上述修正的 Gauss 消去法得到的 $A = LU + R$ 就是关于 \mathcal{F} 和 ω 的松弛不完全 LU 分解.

下面给出修正 Gauss 消去法可以进行下去的充分条件.

定义 6.2 若矩阵 A 满足

$$(1) \ a_{ii} > 0 \ (i = 1, 2, \dots, n-1), \ a_{nn} \geq 0.$$

$$(2) \ a_{ij} \leq 0 \ (i \neq j), \ 1 \leq i, j \leq n.$$

$$(3) \ n(i) = \max\{j : 1 \leq j \leq n, a_{ij} \neq 0\} > i, \ i = 1, 2, \dots, n-1.$$

则称 A 为 \hat{M} 矩阵.

有下面的定理, 详细证明参考文献 [6].

定理 6.4 设 $A = (a_{ij}) \in \mathbb{R}^{n \times n}$ 为弱严格对角占优的 \hat{M} 矩阵, 对任意指标集 \mathcal{F} ($NZ \subset \mathcal{F} \subset ND$) 和 ω ($0 \leq \omega \leq 1$), 修正的 Gauss 消去法给出了 A 关于 \mathcal{F} 和 ω 的松弛不完全 LU 分解.

定义非零模式

$$\tilde{\mathcal{F}} = \mathcal{F} \cup \{(i, i) : i = 1, 2, \dots, n\},$$

它是零模式的补集. 加入修正的 Gauss 消去法是为了确保每步的零模式和非零模式都保持不变. 所以每次修正消元实际上只是对非零模式 $\tilde{\mathcal{F}}$ 的元素进行了 Gauss 消元, 同时将对角元素进行相应的修正. 下面给出松弛不完全 LU 分解具体的算法, 其中 L 存放在 A 的严格下三角部分, U 存放在 A 的上三角部分.

算法 6.5 (松弛不完全 LU 分解) 给定矩阵 A , 非零模式 $\tilde{\mathcal{F}}$ 和松弛参数 ω .

for $k = 1 : n - 1$

 for $i = k + 1 : n$

 if $(i, k) \in \tilde{\mathcal{F}}$ (非零模式)

$$a_{ik} = a_{ik} / a_{kk};$$

 for $j = k + 1 : n$

$$a_{ij} = a_{ij} - a_{ik} \cdot a_{kj};$$

 end

 for $j = k + 1 : n$

 if $(i, j) \notin \tilde{\mathcal{F}}$ (零模式)

$$r_{ij} = r_{ij} + a_{ij}; \ r_{ii} = r_{ii} - \omega * a_{ij};$$

$$a_{ii} = a_{ii} + \omega * a_{ij}; \ a_{ij} = 0;$$

 end

 end

end

end

end

下面给出 ILU 分解的 MATLAB 程序如下:

```

function [L,U,R]=milu(A,NF,omega)
%矩阵A的松弛不完全LU分解
%输入:矩阵A,非零模式NF,松弛参数omega
%输出:单位下三角阵L,上三角阵U,剩余矩阵R,满足A=LU+R
n=size(A,1); R=zeros(n);
for k=1:n-1
    for i=k+1:n
        if (ismember(i+k*sqrt(-1),NF)==1) %(i,k)属于集合NF
            A(i,k)=A(i,k)/A(k,k);
            for j=k+1:n
                A(i,j)=A(i,j)-A(i,k)*A(k,j);
            end
        end
        for j=k+1:n
            if (ismember(i+j*sqrt(-1),NF)==0) %(i,j)不属于集合NF
                R(i,j)=R(i,j)+A(i,j);R(i,i)=R(i,i)-omega*A(i,j);
                A(i,i)=A(i,i)+omega*A(i,j);A(i,j)=0;
            end
        end
    end
end
end
L=tril(A,-1)+eye(n); U=triu(A);

```

由于实际应用中一般不需要剩余矩阵 R 的信息, 为了节省存储量, 可以不对其计算及保存, 这只需在上述算法程序中将涉及计算 R 的有关语句去掉即可.

例 6.7 利用程序 milu.m 对矩阵 A 进行松弛不完全 LU 分解, 其中

$$A = \begin{bmatrix} 4 & 0 & 0 & 1 & 0 \\ 0 & 3 & 2 & 0 & 1 \\ 8 & 0 & 2 & 0 & 0 \\ 0 & 9 & 0 & 5 & 0 \\ 0 & 0 & 2 & 0 & 6 \end{bmatrix},$$

$$\mathcal{F} = \{(1,1), (1,4), (2,2), (2,3), (2,5), (3,1), (3,3), (4,2), (4,4), (5,3), (5,5)\}, \quad \omega = 0.$$

解 编写 MATLAB 脚本文件 ex67.m, 在命令窗口运行之, 可得相应的计算结果(略).

6.3 对称正定方程组的直接法

前面讨论的 Gauss 消去法和 LU 分解法, 都是求解一般方程组的方法, 它们均不考虑方程组系数矩阵本身的特点. 但在实际应用中经常会遇到一些特殊类型的方程组,

其系数矩阵具有某种特殊性, 如对称正定矩阵、稀疏(带状)矩阵等. 对于这些方程组, 若还用原有的一般方法来求解, 势必造成存储空间和计算的浪费. 因此, 有必要构造适合特殊方程组的求解方法. 本节主要介绍解对称正定方程组的 Cholesky 分解法.

6.3.1 Cholesky 分解法

当线性方程组的系数矩阵 A 是对称正定矩阵时, 可利用对称正定的特点使 LU 分解减少计算量, 从而节省存储空间. 由于对称正定矩阵的所有顺序主子式都大于零, 故由定理 6.2 可知 A 存在唯一的 LU 分解. 由于 A 是对称的, 即 $a_{ij} = a_{ji}$, $i, j = 1, 2, \dots, n$. 由 LU 分解式 (6.12)~式 (6.14), 有

$$u_{1i} = a_{1i}, \quad i = 1, \dots, n; \quad l_{i1} = \frac{a_{i1}}{a_{11}}, \quad i = 2, \dots, n,$$

则

$$l_{i1} = \frac{a_{i1}}{a_{11}} = \frac{a_{1i}}{a_{11}} = \frac{u_{1i}}{u_{11}}, \quad i = 2, \dots, n. \quad (6.19)$$

若已求得第 1 步到第 $k-1$ 步的 L 和 U 的元素有如下关系, 即

$$l_{ij} = \frac{u_{ji}}{u_{jj}}, \quad j = 1, \dots, k-1; \quad i = j+1, \dots, n, \quad (6.20)$$

则对于第 k 步, 由式 (6.13)、式 (6.14) 和式 (6.20), 得

$$\begin{aligned} u_{ki} &= a_{ki} - \sum_{r=1}^{k-1} l_{kr} u_{ri} = a_{ki} - \sum_{r=1}^{k-1} \frac{u_{rk} u_{ri}}{u_{rr}}, \quad i = k, \dots, n; \\ l_{ik} &= \left(a_{ik} - \sum_{r=1}^{k-1} l_{ir} u_{rk} \right) / u_{kk} \\ &= \left(a_{ik} - \sum_{r=1}^{k-1} \frac{u_{rk} u_{ri}}{u_{rr}} \right) / u_{kk} = \frac{u_{ki}}{u_{kk}}, \quad i = k+1, \dots, n. \end{aligned}$$

由此, 得

$$l_{ik} = \frac{u_{ki}}{u_{kk}}, \quad k = 1, \dots, n-1; \quad i = k+1, \dots, n. \quad (6.21)$$

这样, 利用式 (6.21) 计算 L 的元素可节省工作量, 计算量节省了将近一半, 而 U 的元素仍用式 (6.13) 计算:

$$u_{ki} = a_{ki} - \sum_{r=1}^{k-1} l_{kr} u_{ri}, \quad k = 2, \dots, n; \quad i = k, \dots, n.$$

注 6.2 由式 (6.21), 得

$$u_{ki} = u_{kk} l_{ik}, \quad k = 1, \dots, n-1; \quad i = k+1, \dots, n,$$

此即

$$U = DL^T,$$

式中: D 为以 $u_{kk} (k=1, \dots, n)$ 为对角元的对角矩阵. 这样, 就把对称正定矩阵 A 分解成了

$$A = LU = LDL^T$$

的形式. 这种分解方法称为 Cholesky 分解法.

下面建立用 Cholesky 分解法求解对称正定方程组的算法步骤.

$$Ax = b \Rightarrow \begin{cases} A = LDL^T, \\ LDL^T x = b, \end{cases} \Rightarrow \begin{cases} Ly = b, \\ Dz = y, \\ L^T x = z. \end{cases} \quad (6.22)$$

算法 6.6 (Cholesky 分解法)

步 1, 输入对称正定矩阵 A 和右端向量 b .

步 2, Cholesky 分解:

$$u_{1i} = a_{1i}, \quad i = 1, \dots, n;$$

$$l_{i1} = u_{1i}/u_{11}, \quad i = 2, \dots, n.$$

对 $k = 2, \dots, n$, 计算

$$u_{ki} = a_{ki} - \sum_{r=1}^{k-1} l_{kr} u_{ri}, \quad i = k, \dots, n;$$

$$l_{ik} = u_{ki}/u_{kk}, \quad i = k+1, \dots, n.$$

步 3, 用向前消去法解下三角方程组 $Ly = b$:

$$y_1 = b_1,$$

$$\text{对 } k = 2, \dots, n, \text{ 计算 } y_k = b_k - \sum_{i=1}^{k-1} l_{ki} y_i.$$

步 4, 解对角形方程组 $Dz = y$:

$$\text{对 } k = 1, \dots, n, \text{ 计算 } z_k = y_k/d_k.$$

步 5, 用回代法解上三角方程组 $L^T x = z$:

$$x_n = z_n,$$

$$\text{对 } k = n-1, \dots, 1, \text{ 计算 } x_k = z_k - \sum_{i=k+1}^n l_{ik} x_i.$$

根据算法 6.6, 编制 MATLAB 程序如下:

```
%Cholesky分解法程序-mschol.m
function [x,L,D]=mschol(A,b)
%用Cholesky分解法解对称正定方程组Ax=b
%输入:系数矩阵A,右端项b
%输出:解向量x,单位下三角阵L,对角阵D
n=size(A,1); D=zeros(1,n); L=eye(n,n);
U(1,:)=A(1,:); L(2:n,1)=U(1,2:n)/U(1,1); %Cholesky分解
for k=2:n
    U(k,k:n)=A(k,k:n)-L(k,1:k-1)*U(1:k-1,k:n);
```



```

    L(k+1:n,k)=U(k,k+1:n)/U(k,k);
end
%求解下三角方程组Ly=b(向前消去法)
y=zeros(n,1); y(1)=b(1);
for k=2:n,
    y(k)=b(k)-L(k,1:k-1)*y([1:k-1]);
end
%求解对角方程组Dz=y
D=diag(diag(U));
for k=1:n,
    z(k)=y(k)/D(k,k);
end
%求解上三角方程组L'x=z(回代法)
x=zeros(n,1); U=L'; x(n)=z(n);
for k=(n-1):-1:1,
    x(k)=z(k)-U(k,k+1:n)*x(k+1:n);
end

```

例 6.8 利用程序 mschol.m 计算下列对称正定方程组的解

$$\begin{bmatrix} 1 & 1 & -1 \\ 1 & 2 & -3 \\ -1 & -3 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ -3 \\ 2 \end{bmatrix}.$$

解 在 MATLAB 命令窗口输入:

```

>> A=[1,1,-1; 1,2,-3; -1,-3,3];
>> b=[0,-3,2]';
>> [x,L,D]=mschol(A,b)

```

可得计算结果 (略).

6.3.2 不完全 Cholesky 分解

在预处理共轭梯度法中, 往往需要对称正定矩阵 A 的不完全 Cholesky 分解 $A \approx LDL^T$ 来获得预处理子 $M = LDL^T$. 由于 A 的对称性, 此时, 假定指标集 \mathcal{F} 也是对称正定的, 即若 $(i, j) \in \mathcal{F}$, 则必有 $(j, i) \in \mathcal{F}$. 相对于非对称情形的松弛 ILU 分解, 可以考虑 A 关于 \mathcal{F} 和 ω 的松弛不完全 Cholesky 分解 (简称 RIC 分解):

$$A = LDL^T + R, \quad (6.23)$$

式中: L 为单位下三角矩阵且满足式 (6.15); R 为剩余矩阵并满足式 (6.17) 和式 (6.18); D 为对角矩阵.

易知, 此时 R 也是对称的, 而且也可以采用前面介绍的修正 Gauss 消去法实现分解 (6.23), 即先用修正 Gauss 消去法求 A 的 RILU 分解

$$A = LU + R,$$

然后取 $D = \text{diag}(u_{11}, u_{22}, \dots, u_{nn})$ 为 U 的对角元素所构成的对角矩阵. 这样便得到了式 (6.23) 中的 L , D 和 R .

下面的定理给出了对称正定矩阵 A 关于 \mathcal{F} 和 ω 的松弛不完全 Cholesky 分解 (6.23) 中 LDL^T 正定的一个充分条件.

可以在算法 6.5 的基础上, 利用 A 的对称性给出求 A 的 RIC 分解的算法. 也可从 Cholesky 分解出发稍加修正求得. 下面是不完全 Cholesky 分解的一个简易可行的算法.

算法 6.7 (不完全 Cholesky 分解法) 给定对称正定矩阵 A . 下面的算法计算 A 的不完全 Cholesky 分解: $A \approx LL^T$.

```

n = size(A, 1);
for k = 1 : n
    lkk = (akk - ∑r=1k-1 lkr2)1/2;
    for i = k + 1 : n
        if (aik = 0)
            lik = 0;
        else
            lik = (aik - ∑r=1k-1 lirlkr) / lkk;
        end
    end
end
end

```

显然, 由算法 6.7 得到的不完全分解 $M = LL^T$ 与 A 有相同的稀疏性, 但这一算法并不总是稳定的 (当 l_{kk} 很小的时候, $k = 1, 2, \dots, n$).

算法 6.7 的 MATLAB 程序如下:

```

%不完全Cholesky分解程序-michol.m
function [L]=michol(A)
%输入:对称正定矩阵A
%输出:下三角矩阵L,满足A=LL'+R
n=size(A,1); L=zeros(n,n);
L(1,1)=sqrt(A(1,1)); L(2,1)=A(2,1)/L(1,1);
for k=2:n
    s=0;
    for(p=1:k-1),s=s+L(k,p)^2;end
    L(k,k)=sqrt(A(k,k)-s);

```

```

for i=k+1:n
    if (A(i,k)~=0)
        s=0;
        for (p=1:k-1),s=s+L(i,p)*L(k,p); end
        L(i,k)=(A(i,k)-s)/L(k,k);
    end
end
end
end

```

例 6.9 利用程序 michol.m 对矩阵 A 进行松弛不完全 Cholesky 分解, 其中

$$A = \begin{bmatrix} 8 & -2 & -1 & 0 & 0 & 1 \\ -2 & 8 & -2 & 0 & 3 & 0 \\ -1 & -2 & 8 & 0 & 0 & 0 \\ 0 & 0 & 0 & 8 & -2 & -1 \\ 0 & 3 & 0 & -2 & 8 & -2 \\ 1 & 0 & 0 & -1 & -2 & 8 \end{bmatrix}.$$

解 编写 MATLAB 脚本文件 ex69.m, 在命令窗口运行之, 可得相应的计算结果 (略).

6.4 带状线性方程组的直接法

6.4.1 三对角方程组

在科学与工程计算中, 经常遇到求解三对角方程组的问题. 例如, 三次样条插值计算, 以及用有限差分法求解二阶常系数线性常微分方程的边值问题和热传导问题, 经常需要求解三对角方程组 (即系数矩阵为三对角矩阵的方程组). 三对角矩阵属于所谓的“带状矩阵”, 在大多数应用中, 带状矩阵是严格对角占优的或正定的. 下面给出带状矩阵的定义.

定义 6.3 n 阶矩阵称为带状矩阵, 如果存在正整数 p, q ($1 < p, q < n$), 当 $j \geq i+p$ 或 $i \geq j+q$ 时, 有 $a_{ij} = 0$, 并称 $w = p+q-1$ 为该带状矩阵的“带宽”. n 阶矩阵称为带状矩阵, 如果存在正整数 p, q ($1 < p, q < n$), 当 $i+p \leq j$ 或 $j+q \leq i$ 时, 有 $a_{ij} = 0$, 并称 $w = p+q-1$ 为该带状矩阵的“带宽”.

三对角方程组的一般形式是

$$Ax := \begin{bmatrix} b_1 & c_1 & & & \\ a_2 & b_2 & c_2 & & \\ & \ddots & \ddots & \ddots & \\ & & a_{n-1} & b_{n-1} & c_{n-1} \\ & & & a_n & b_n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n-1} \\ x_n \end{bmatrix} = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_{n-1} \\ f_n \end{bmatrix} := f. \quad (6.24)$$

显然三对角矩阵的带宽为 3. 本节介绍求解方程组 (6.24) 的追赶法和变参数追赶法.

1. 追赶法

将顺序 LU 分解法应用于三对角方程组得到所谓的“追赶法”. 事实上, 一方面, 可将三对角矩阵 A 分解为

$$A = LU, \quad (6.25)$$

式中:

$$L = \begin{bmatrix} l_1 & & & & \\ a_2 & l_2 & & & \\ & \ddots & \ddots & & \\ & & a_{n-1} & l_{n-1} & \\ & & & a_n & l_n \end{bmatrix}, \quad U = \begin{bmatrix} 1 & u_1 & & & \\ & 1 & u_2 & & \\ & & \ddots & \ddots & \\ & & & 1 & u_{n-1} \\ & & & & 1 \end{bmatrix}.$$

比较式 (6.25) 两端的对应元素, 得

$$b_1 = l_1, \quad c_{k-1} = l_{k-1}u_{k-1}, \quad b_k = a_k u_{k-1} + l_k, \quad k = 2, \dots, n.$$

于是有

$$l_1 = b_1, \quad u_{k-1} = \frac{c_{k-1}}{l_{k-1}}, \quad l_k = b_k - a_k u_{k-1}, \quad k = 2, \dots, n. \quad (6.26)$$

另一方面, 方程组 (6.24) 等价于求解 $Ly = f$ 和 $Ux = y$, 其中 $y = (y_1, y_2, \dots, y_n)^T$. 分别比较 $Ly = f$ 和 $Ux = y$ 两端的对应元素, 得

$$\begin{cases} l_1 y_1 = f_1, & a_k y_{k-1} + l_k y_k = f_k, & k = 2, 3, \dots, n, \\ x_k + u_k x_{k+1} = y_k, & k = 1, 2, \dots, n-1, & x_n = y_n. \end{cases} \quad (6.27)$$

结合式 (6.26) 和式 (6.27) 可得下面的算法.

算法 6.8 (追赶法) 计算三对角方程组 $Ax = f$ 的解.

```

 $l_1 = b_1; \quad y_1 = \frac{f_1}{l_1};$ 
for  $k = 2 : n$ 
     $u_{k-1} = \frac{c_{k-1}}{l_{k-1}};$ 
     $l_k = b_k - a_k u_{k-1};$ 
     $y_k = \frac{f_k - a_k y_{k-1}}{l_k};$ 
end
 $x_n = y_n;$ 
for  $k = n - 1 : -1 : 1$ 
     $x_k = y_k - u_k x_{k+1};$ 
end

```

算法 6.8 被称为“追赶法”的原因: 一是第 1 步关于指标 k 由小到大计算 l_k , u_k 和 y_k 这三个量, 这是“向前追”的过程; 二是第 2 步关于指标 k 从大到小计算方程组的解 x_k , 此即“往回赶”的过程. 追赶法只有 $5n-4$ 次乘除法运算和 $3n-3$ 次加减法运算, 且当系数矩阵对角占优时数值稳定, 是解三对角方程组的优秀算法. 编程计算时, 可将 l_k , u_k 依次存放在 b_k , c_k 的位置, 而将 y_k 和 x_k 先后存放在 f_k 的位置, 因此整个计算过程只需 $4n$ 个存储单元.

追赶法的 MATLAB 程序如下:

```
%追赶法程序-mchase.m
function [f]=mchase(a,b,c,f)
%用追赶法解三对角方程组Ax=f
%输入:a为A的下对角线,b为A的主对角线,c为A的次上对角线,f为右端向量
%输出:解向量f(LU分解中的l(k),u(k)存放在b(k),c(k)的位置,
%      y(k)和x(k)先后存放在f(k)的位置)
n=length(b); f(1)=f(1)/b(1);
for k=2:n
    c(k-1)=c(k-1)/b(k-1);
    b(k)=b(k)-a(k)*c(k-1);
    f(k)=(f(k)-a(k)*f(k-1))/b(k);
end
for k=n-1:-1:1
    f(k)=f(k)-c(k)*f(k+1);
end
```

定理 6.5 若方程组 (6.24) 的系数矩阵的元素满足条件

$$|b_1| > |c_1| > 0, |b_n| > |a_n| > 0, |b_i| > |a_i| + |c_i|, \quad i = 2, \dots, n-1,$$

则追赶法是可行的.

证明 由式 (6.26) 和式 (6.27) 可知, 只需证明 $l_k \neq 0$ ($k = 1, 2, \dots, n$) 即可. 显然 $l_1 = b_1 \neq 0$, 且 $|l_1| = |b_1| > |c_1|$. 设 $|l_{k-1}| > |c_{k-1}|$, 则

$$\begin{aligned} |l_k| &= |b_k - a_k u_{k-1}| = \left| b_k - \frac{a_k}{l_{k-1}} c_{k-1} \right| \\ &\geq |b_k| - |a_k| \cdot \left| \frac{c_{k-1}}{l_{k-1}} \right| > |b_k| - |a_k| \\ &> \begin{cases} |c_k|, & k < n, \\ 0, & k = n, \end{cases} \end{aligned}$$

即 $l_k \neq 0$, $k = 2, \dots, n$. 从而, 追赶法是可行的. 证毕. \square

注 6.3 满足定理 6.5 条件的三对角矩阵即为严格对角占优的, 例 6.5 说明对于严格对角占优的三对角矩阵, 追赶法总是可行的.

2. 变参数追赶法

需要指出的是, 追赶法的本质是没有选取主元的 Gauss 消去法, 因此对于一般的三对角矩阵, 不一定存在如式 (6.25) 的三角分解. 或者三角分解可以进行下去, 但计算过程不是数值稳定的, 因而最终得到的解并不可靠. 这就促使考虑对追赶法进行改进和修正.

将三对角矩阵 A 分解为

$$A = D\tilde{L}\tilde{U}, \quad (6.28)$$

式中: $D = \text{diag}(d_1, d_2, \dots, d_n)$,

$$\tilde{L} = \begin{bmatrix} l_1 & 1 & & & \\ & l_2 & 1 & & \\ & & \ddots & \ddots & \\ & & & l_n & 1 \end{bmatrix}_{n \times (n+1)}, \quad \tilde{U} = \begin{bmatrix} u_1 & & & & \\ 1 & u_2 & & & \\ & 1 & \ddots & & \\ & & \ddots & u_n & \\ & & & 1 \end{bmatrix}_{(n+1) \times n}$$

比较式 (6.25) 两端的元素, 得

$$\begin{cases} b_k = d_k(l_k u_k + 1), & k = 1, 2, \dots, n; \\ a_k = d_k l_k, & c_{k-1} = d_{k-1} u_k, & k = 2, \dots, n. \end{cases}$$

选取 l_1, u_1 使得 $l_1 u_1 + 1 \neq 0$, 则有

$$\begin{cases} d_1 = \frac{b_1}{l_1 u_1 + 1}, \\ u_k = \frac{c_{k-1}}{d_{k-1}}, & d_k = b_k - a_k u_k, & l_k = \frac{a_k}{d_k}, \\ & (k = 2, 3, \dots, n). \end{cases}$$

注意到, 求解 $Ax = f$ 等价于求解 $\tilde{L}y = D^{-1}f := g$ 和 $\tilde{U}x = y$, 其中

$$y = (y_1, y_2, \dots, y_{n+1})^T, \quad g = \left(\frac{f_1}{d_1}, \frac{f_2}{d_2}, \dots, \frac{f_n}{d_n} \right)^T.$$

比较 $\tilde{L}y = g$ 和 $\tilde{U}x = y$ 两端的元素, 得

$$l_k y_k + y_{k+1} = g_k, \quad k = 1, 2, \dots, n; \quad (6.29)$$

$$\begin{cases} u_1 x_1 = y_1, \\ x_k + u_{k+1} x_{k+1} = y_{k+1}, & k = 1, 2, \dots, n-1, \\ x_n = y_{n+1}. \end{cases} \quad (6.30)$$

由式 (6.29) 和式 (6.30) 可以发现, 只要计算出 x_1 , 其他的可以依次计算出来. 下面推导计算 x_1 的公式.

由式 (6.30), 得

$$\begin{aligned}
 x_1 &= y_2 - u_2 x_2 = y_2 - u_2(y_3 - u_3 x_3) \\
 &= y_2 + (-u_2)y_3 + (-u_2)(-u_3)x_3 \\
 &= \cdots \\
 &= y_2 + (-u_2)y_3 + (-u_2)(-u_3)y_4 \\
 &\quad + \cdots + (-u_2)(-u_3) \cdots (-u_{n-1})y_n + (-u_2)(-u_3) \cdots (-u_n)y_{n+1}.
 \end{aligned} \tag{6.31}$$

由式 (6.29), 得

$$\begin{aligned}
 y_{k+1} &= g_k - l_k y_k = g_k - l_k(g_{k-1} + l_{k-1}y_{k-1}) \\
 &= g_k + (-l_k)g_{k-1} + (-l_k)(-l_{k-1})y_{k-1} \\
 &= \cdots \\
 &= g_k + (-l_k)g_{k-1} + (-l_k)(-l_{k-1})g_{k-2} \\
 &\quad + \cdots + (-l_k)(-l_{k-1}) \cdots (-l_2)g_1 \\
 &\quad + (-l_k)(-l_{k-1}) \cdots (-l_1)(u_1 x_1), \quad k = 1, 2, \cdots, n.
 \end{aligned} \tag{6.32}$$

将式 (6.32) 代入式 (6.31) 并整理, 得

$$\begin{aligned}
 &[1 + u_1 l_1 + (u_1 u_2)(l_2 l_1) + \cdots + (u_1 \cdots u_n)(l_n \cdots l_1)]x_1 \\
 &= [1 + u_2 l_2 + \cdots + (u_2 \cdots u_n)(l_n \cdots l_2)]g_1 \\
 &\quad + (-u_2)[1 + u_3 l_3 + \cdots + (u_3 \cdots u_n)(l_n \cdots l_3)]g_2 \\
 &\quad + \cdots + [(-u_2) \cdots (-u_{n-1})](1 + u_n l_n)g_{n-1} + [(-u_2) \cdots (-u_n)]g_n.
 \end{aligned}$$

令

$$\begin{aligned}
 s_k &= 1 + u_k l_k + \cdots + (u_k \cdots u_n)(l_n \cdots l_k), \quad k = 1, 2, \cdots, n, \\
 t_1 &= s_2 g_1 + (-u_2)s_3 g_2 + \cdots + [(-u_2) \cdots (-u_{n-1})]s_n g_{n-1} + [(-u_2) \cdots (-u_n)]g_n.
 \end{aligned}$$

则

$$x_1 = \frac{t_1}{s_1},$$

且有递推公式

$$\begin{aligned}
 s_n &= 1 + u_n l_n, \quad s_k = 1 + u_k s_{k+1} l_k, \quad k = n-1, n-2, \cdots, 1, \\
 t_n &= g_n, \quad t_k = s_{k+1} g_k - u_{k+1} t_{k+1}, \quad k = n-1, n-2, \cdots, 1.
 \end{aligned}$$

综合上述, 可得下面的变参数追赶法.

算法 6.9 (变参数追赶法) 计算三对角方程组 $Ax = f$ 的解. 选取 l_1 和 u_1 使得 $l_1 u_1 + 1 \neq 0$.

```


$$d_1 = \frac{b_1}{l_1 u_1 + 1}; g_1 = \frac{f_1}{d_1};$$

for  $k = 2 : n$ 
    
$$u_k = \frac{c_{k-1}}{d_{k-1}}; d_k = b_k - a_k u_k; l_k = \frac{a_k}{d_k}; g_k = \frac{f_k}{d_k};$$

end

 $s_n = 1 + u_n l_n; t_n = g_n;$ 
for  $k = n - 1 : -1 : 1$ 
    
$$s_k = 1 + u_k s_{k+1} l_k; t_k = s_{k+1} g_k - u_{k+1} t_{k+1};$$

end


$$x_1 = \frac{t_1}{s_1}; y_1 = u_1 x_1;$$

for  $k = 1 : n$ 
    
$$y_{k+1} = g_k - l_k y_k;$$

end

 $x_n = y_{n+1};$ 
for  $k = n - 1 : -1 : 2$ 
    
$$x_k = y_{k+1} - u_{k+1} x_{k+1};$$

end

```

算法 6.9 需 $10n - 8$ 次乘除运算, $5n - 5$ 次加减运算. 由于算法中的 l_1 和 u_1 的选取只需使得 $l_1 u_1 + 1 \neq 0$, 可以认为它们是两个可变的参数, 故称为“变参数追赶法”.

变参数追赶法的 MATLAB 程序如下:

```

%变参数追赶法程序-mchase_var.m
function [x]=mchase_var(a,b,c,f,l1,u1)
%用变参数追赶法解三对角方程组Ax=f
%输入:a为A的下对角线,b为A的主对角线,c为A的次
%    上对角线,f为右端向量;l1,u1满足l1*u1+1不为0.
%输出:解向量x
n=length(b);x=zeros(n,1);l=zeros(n,1);
u=zeros(n,1);l(1)=l1; u(1)=u1;
d(1)=b(1)/(l(1)*u(1)+1); g(1)=f(1)/d(1);
for k=2:n
    u(k)=c(k-1)/d(k-1);d(k)=b(k)-a(k)*u(k);
    l(k)=a(k)/d(k);g(k)=f(k)/d(k);
end
s(n)=1+u(n)*l(n); t(n)=g(n);
for k=n-1:-1:1
    s(k)=1+u(k)*s(k+1)*l(k);
    t(k)=s(k+1)*g(k)-u(k+1)*t(k+1);

```



```
end
x(1)=t(1)/s(1);y(1)=u(1)*x(1);
for k=1:n
    y(k+1)=g(k)-l(k)*y(k);
end
x(n)=y(n+1);
for k=n-1:-1:2
    x(k)=y(k+1)-u(k+1)*x(k+1);
end
```

例 6.10 设

$$A = \begin{bmatrix} 4 & -1 & & & \\ -2 & 4 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -2 & 4 & -1 \\ & & & -2 & 4 \end{bmatrix}, \quad f = \begin{bmatrix} 3 \\ 1 \\ \vdots \\ 1 \\ 2 \end{bmatrix}.$$

用追赶法和变参数追赶法 (取 $l_1 = u_1 = 1$) 求解三对角方程组 $Ax = f$, 并计算实际误差 $\|f - Ax\|_2$.

解 利用程序 mchase.m 和 mchase_var.m, 编写 MATLAB 脚本文件 ex610.m, 并对不同的维数 n , 在命令窗口运行该程序, 得到数值结果如表 6.1 所示.

表 6.1 追赶法和变参数追赶法的数值结果

维数	追赶法		变参数追赶法	
	误差	计算时间	误差	计算时间
1024	7.1650e-15	0.0056	1.2212e-15	0.0103
2048	1.0091e-14	0.0058	1.2212e-15	0.0106
4096	1.4241e-14	0.0064	1.2212e-15	0.0109
8192	2.0118e-14	0.0066	1.2212e-15	0.0128

例 6.11 设

$$A = \begin{bmatrix} 2 & 3 & & & \\ 6 & 9 & 3 & & \\ & \ddots & \ddots & \ddots & \\ & & 6 & 9 & 3 \\ & & & 6 & 9 \end{bmatrix}, \quad f = \begin{bmatrix} 5 \\ 10 \\ \vdots \\ 10 \\ 12 \end{bmatrix}.$$

试用追赶法和变参数追赶法 (取 $l_1 = u_1 = 1$) 求解三对角方程组 $Ax = f$, 并计算实际误差 $\|f - Ax\|_2$.

解 对于此例, 追赶法已经失效, 但变参数追赶法仍然非常有效, 数值结果 (运行 MATLAB 脚本文件 ex611.m) 如表 6.2 所示.

表 6.2 变参数追赶法的数值结果

维数	追赶法	变参数追赶法	
		误差	计算时间
1024	失效	5.7293e-14	0.0101
4096	失效	1.1391e-13	0.0111

6.4.2 块三对角方程组

用有限差分法五点格式离散二维 Poisson 方程, 得到的代数方程组 $Ax = f$ 的系数矩阵是一个块三对角矩阵, 即

$$\begin{bmatrix} B_1 & C_1 & & & \\ A_2 & B_2 & C_2 & & \\ & \ddots & \ddots & \ddots & \\ & & A_{m-1} & B_{m-1} & C_{m-1} \\ & & & A_m & B_m \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{m-1} \\ x_m \end{bmatrix} = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_{m-1} \\ f_m \end{bmatrix}, \quad (6.33)$$

式中: A_k, B_k, C_k 均为 r 阶方阵; x_k, f_k 为 r 维列向量. 当 A 可逆时, 方程组有唯一解 x^* . 本节介绍求解方程组 (6.33) 的块追赶法和双参数法.

1. 块追赶法

将三对角方程矩阵 A 分解为 $A = LU$, 其中

$$L = \begin{bmatrix} L_1 & & & & \\ A_2 & L_2 & & & \\ & \ddots & \ddots & \ddots & \\ & & A_{m-1} & L_{m-1} & \\ & & & A_m & L_m \end{bmatrix}, \quad U = \begin{bmatrix} I & U_1 & & & \\ & I & U_2 & & \\ & & \ddots & \ddots & \\ & & & I & U_{m-1} \\ & & & & I \end{bmatrix}.$$

比较 $A = LU$ 两端的对应的子矩阵, 得

$$B_1 = L_1, \quad C_{k-1} = L_{k-1}U_{k-1}, \quad B_k = A_kU_{k-1} + L_k, \quad k = 2, \dots, n.$$

于是有

$$L_1 = B_1, \quad U_{k-1} = L_{k-1}^{-1}C_{k-1}, \quad L_k = B_k - A_kU_{k-1}, \quad k = 2, \dots, n. \quad (6.34)$$

求解三对角方程组 $Ax = f$ 等价于求解

$$Ly = f \quad \text{和} \quad Ux = y,$$

式中: $y = [y_1^T, y_2^T, \dots, y_m^T]^T$. 分别比较 $Ly = f$ 和 $Ux = y$ 两端的对应元素, 得

$$\begin{cases} L_1 y_1 = f_1, & A_k y_{k-1} + L_k y_k = f_k, & k = 2, 3, \dots, m, \\ x_k + U_k x_{k+1} = y_k, & k = 1, 2, \dots, m-1, & x_m = y_m. \end{cases} \quad (6.35)$$

结合式 (6.34) 和式 (6.35) 可得下面的算法.

算法 6.10 (块追赶法) 计算块三对角方程组 $Ax = f$ 的解.

$L_1 = B_1; y_1 = L_1^{-1} f_1;$

for $k = 2 : m$

$U_{k-1} = L_{k-1}^{-1} C_{k-1}; L_k = B_k - A_k U_{k-1}; y_k = L_k^{-1} (f_k - A_k y_{k-1});$

end

$x_m = y_m;$

for $k = m-1 : -1 : 1$

$x_k = y_k - U_k x_{k+1};$

end

使用块追赶法求解三对角方程组 $Ax = f$ 需要的乘除法次数约为

$$mr^2 \left(\frac{10}{3}r + 3 \right) - r^2(3r + 2).$$

编程计算时, 可将 L_k 和 U_k 依次存放在 B_k 和 C_k 的位置, 而将 y_k 和 x_k 先后存放在 f_k 的位置.

块追赶法的 MATLAB 程序如下:

```
%块追赶法程序-mchase_block.m
function [fi]=mchase_block(Ai,Bi,Ci,fi,m)
%用块追赶法解块三对角方程组Ax=f
%输入: Ai为A的次下对角块, Bi为A的主对角块,
%      Ci为A的次上对角块, fi为右端向量, m为分块数
%输出: 解向量fi, 其中LU分解的L{k}, U{k}存放在Bi{k},
%      Ci{k}的位置, y{k}和x{k}先后存放在fi{k}的位置
fi{1}=Bi{1}\fi{1};
for k=2:m
    Ci{k-1}=Bi{k-1}\Ci{k-1}; Bi{k}=Bi{k}-Ai{k}*Ci{k-1};
    fi{k}=Bi{k}\(fi{k}-Ai{k}*fi{k-1});
end
for k=m-1:-1:1
    fi{k}=fi{k}-Ci{k}*fi{k+1};
end
```

2. 双参数法

注意到块三角对方程组 (6.33) 的前 $m-1$ 个子方程为

$$\begin{cases} B_1x_1 + C_1x_2 = f_1, \\ A_2x_1 + B_2x_2 + C_2x_3 = f_2, \\ \vdots \\ A_{m-1}x_{m-2} + B_{m-1}x_{m-1} + C_{m-1}x_m = f_{m-1}. \end{cases}$$

将上式中的 x_2, x_3, \dots, x_m 都用 x_1 表示出来, 即

$$\begin{aligned} x_2 &= C_1^{-1}(f_1 - B_1x_1) := s_2 + T_2x_1, \\ x_3 &= C_2^{-1}(f_2 - A_2x_1 - B_2x_2) \\ &= C_2^{-1}(f_2 - B_2s_2) - C_2^{-1}(A_2 + B_2T_2)x_1 := s_3 + T_3x_1, \\ &\vdots \\ x_m &= C_{m-1}^{-1}(f_{m-1} - A_{m-1}x_{m-2} - B_{m-1}x_{m-1}) := s_m + T_mx_1. \end{aligned}$$

记 $s_1 = 0, T_1 = I_r$ (r 阶单位阵), 则 $x_1 = s_1 + T_1x_1$. 于是形式上有

$$x_k = s_k + T_kx_1, \quad k = 1, 2, \dots, m. \quad (6.36)$$

下面推导参数序列 $\{s_k\}$ 和 $\{T_k\}$ 的递推计算公式. 注意到由第 k 个子方程 $A_kx_{k-1} + B_kx_k + C_kx_{k+1} = f_k$ 可得

$$\begin{aligned} x_{k+1} &= C_k^{-1}(f_k - A_kx_{k-1} - B_kx_k) \\ &= C_k^{-1}(f_k - A_ks_{k-1} - B_ks_k) - C_k^{-1}(A_kT_{k-1} + B_kT_k)x_1. \end{aligned}$$

将上式与式 (6.36) 比较, 得

$$\begin{aligned} s_1 &= 0, \quad s_2 = C_1^{-1}f_1, \\ s_{k+1} &= C_k^{-1}(f_k - A_ks_{k-1} - B_ks_k), \quad k = 2, \dots, m-1, \\ T_1 &= I_r, \quad T_2 = -C_1^{-1}B_1, \\ T_{k+1} &= -C_k^{-1}(A_kT_{k-1} + B_kT_k), \quad k = 2, \dots, m-1. \end{aligned}$$

将式 (6.36) 代入式 (6.33) 的第 m 个子方程, 得

$$(A_mT_{m-1} + B_mT_m)x_1 = f_m - A_ms_{m-1} - B_ms_m. \quad (6.37)$$

若 r 阶方程组 (6.37) 有解 x_1 , 则可由式 (6.36) 确定方程组 $Ax = f$ 的唯一解. 称这种方法为解块三对角方程组的双参数法.

因为 A 可逆, 故方程组 $Ax = f$ 有唯一解, 从而 x_1 存在. 假设方程组 (6.37) 有无穷多组解, 则对每一个解 x_1 , 由式 (6.36) 都可求得方程组 $Ax = f$ 的一组解, 这与

方程组 $Ax = f$ 有唯一解矛盾. 因此 x_1 存在且唯一, 从而方程组 (6.37) 的系数矩阵 $A_m T_{m-1} + B_m T_m$ 非奇异, 且有

$$x_1 = (A_m T_{m-1} + B_m T_m)^{-1} (f_m - A_m s_{m-1} - B_m s_m),$$

代入 (6.36) 即可求得方程组 $Ax = f$ 的唯一解. 算法如下.

算法 6.11 (双参数法) 计算块三对角方程组 $Ax = f$ 的解.

$s_1 = 0; s_2 = C_1^{-1} f_1; T_1 = I_r; T_2 = -C_1^{-1} B_1;$

for $k = 2 : m - 1$

$s_{k+1} = C_k^{-1} (f_k - A_k s_{k-1} - B_k s_k);$

$T_{k+1} = -C_k^{-1} (A_k T_{k-1} + B_k T_k);$

end

$x_1 = (A_m T_{m-1} + B_m T_m)^{-1} (f_m - A_m s_{m-1} - B_m s_m);$

for $k = 2 : m$

$x_k = s_k + T_k x_1;$

end

使用双参数法求解三对角方程组 $Ax = f$ 需要的乘除法次数约为

$$mr^2 \left(\frac{13}{3} r + 4 \right) - r^2 (4r + 3).$$

双参数法的 MATLAB 程序如下:

```
%双参数法程序-mdouble_par.m
function [x]=mdouble_par(Ai,Bi,Ci,fi,m)
%用双参数法解三对角方程组Ax=f
%输入: Ai为A的次下对角块, Bi为A的主对角块,
%      Ci为A的次上对角块, fi为右端向量
%输出: 解向量x
x=cell(m,1); s=cell(m,1); T=cell(m,1);
s{1}=zeros(3,1); s{2}=Ci{1}\fi{1};
T{1}=eye(3); T{2}=-Ci{1}\Bi{1};
for k=2:m-1
    s{k+1}=Ci{k}\(fi{k}-Ai{k}*s{k-1}-Bi{k}*s{k});
    T{k+1}=-Ci{k}\(Ai{k}*T{k-1}+Bi{k}*T{k});
end
x{1}=(Ai{m}*T{m-1}+Bi{m}*T{m})\(fi{m}-Ai{m}*s{m-1}-Bi{m}*s{m});
for k=2:m
    x{k}=s{k}+T{k}*x{1};
end
```

例 6.12 选取子矩阵分别为

$$B_k = \begin{bmatrix} 4 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 4 \end{bmatrix}, \quad A_k = C_k^T = \begin{bmatrix} 13 & 0 & 0 \\ 0 & 11 & 0 \\ 1 & 0 & 12 \end{bmatrix}, \quad k = 1, 2, \dots, m.$$

分别对 $m = 10^3, 5 \times 10^3, 10^4, 5 \times 10^4, 10^5, 5 \times 10^5$ 使用块追赶法和双参数法求解块三对角矩阵 $Ax = f$, 其中 $f_k = (1, 0, 1)^T$. 并将所得数值解 x 代入 $Ax = f$, 各子方程的误差为

$$\delta_1 = \|B_1 x_1 + C_1 x_2 - f_1\|_2,$$

$$\delta_k = \|A_k x_{k-1} + B_k x_k + C_k x_{k+1} - f_k\|_2, \quad k = 2, \dots, m-1,$$

$$\delta_m = \|A_m x_{m-1} + B_m x_m - f_m\|_2.$$

两种方法的计算时间 (PC Intel Core i5-3470 CPU 3.2GHz, MATLAB R2015b) 和各子方程的误差最大值 $\max\{\delta_k\}_{k=1}^m$ 如表 6.3 所示.

表 6.3 块追赶法和双参数法的数值结果

维数	块追赶法		双参数法	
	误差	计算时间	误差	计算时间
1×10^3	1.9956e-11	0.0319	5.5943e-16	0.0235
5×10^3	1.8646e-11	0.1277	7.0217e-16	0.0821
1×10^4	1.2669e-10	0.2503	4.9772e-16	0.1551
5×10^4	5.1618e-11	1.2301	8.8991e-16	0.7532
1×10^5	1.2903e-09	2.4382	8.9509e-16	1.4845
5×10^5	8.6523e-10	12.1268	6.2942e-16	7.4758

从表 6.3 可以看出, 对于此例, 双参数法要比块追赶法有效得多.

例 6.13 选取子矩阵分别为

$$B_k = \begin{bmatrix} 2 & -1 & 0 \\ -2 & 1 & 0 \\ 0 & 0 & 3 \end{bmatrix}, \quad A_k = C_k = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{bmatrix}, \quad k = 1, 2, \dots, m.$$

分别对 $m = 10^3, 5 \times 10^3, 10^4, 5 \times 10^4, 10^5, 5 \times 10^5$ 使用块追赶法和双参数法求解块三对角矩阵 $Ax = f$, 其中 $f_k = (1, 2, 1)^T$. 两种方法的计算时间 (PC Intel Core i5-3470 CPU 3.2GHz, MATLAB R2015b) 和各子方程的误差最大值 $\max\{\delta_k\}_{k=1}^m$ 如表 6.4 所示.

从表 6.4 可以看出, 当块追赶法失效时, 双参数法依然十分有效. 当然, 也可以构造出块追赶法有效而双参数法失效的算例. 由此可见, 块追赶法和双参数法可以互为补充.

表 6.4 双参数法的数值结果

维数	块追赶法	双参数法	
		误差	计算时间
1×10^3	失效	4.4409e-16	0.0240
5×10^3	失效	4.4409e-16	0.0848
1×10^4	失效	5.5511e-16	0.1622
5×10^4	失效	4.4409e-16	0.7594
1×10^5	失效	6.6613e-16	1.5165
5×10^5	失效	5.5511e-16	7.4960

6.5 直接法的舍入误差分析

本节对用直接法求解方程组得到的解进行舍入误差分析. 下面先介绍矩阵的条件数, 它是判断矩阵病态与否的一种度量.

6.5.1 矩阵的条件数

定义 6.4 设 A 为非奇异矩阵, 称 $\kappa(A) = \|A^{-1}\| \cdot \|A\|$ 为矩阵 A 的条件数, 这里 $\|\cdot\|$ 是任意的算子范数.

从定义 6.4 知道, $\kappa(A) = \kappa(A^{-1})$. 常用的条件数有

(1) 无穷范数条件数

$$\kappa(A)_\infty = \|A^{-1}\|_\infty \cdot \|A\|_\infty.$$

(2) 谱条件数

$$\kappa(A)_2 = \|A^{-1}\|_2 \cdot \|A\|_2 = \sqrt{\frac{\lambda_{\max}(A^T A)}{\lambda_{\min}(A^T A)}}.$$

当 A 为对称矩阵时, $\kappa(A)_2 = |\lambda_1|/|\lambda_n|$, 其中 λ_1 和 λ_n 分别是 A 的绝对值最大和绝对值最小的特征值.

容易验证, 矩阵的条件数有如下性质:

- (1) 任意非零矩阵的条件数 $\kappa(A) \geq 1$.
- (2) 若 $c \neq 0$ 为常数, 则 $\kappa(cA) = \kappa(A)$.
- (3) 若 A 为正交矩阵, 则 $\kappa(A)_2 = 1$.
- (4) 若 A 为非奇异矩阵, Q 为正交矩阵, 则

$$\kappa(QA)_2 = \kappa(AQ)_2 = \kappa(A)_2.$$

6.5.2 矩阵条件数的估算

本节介绍一种非常实用的方法来估计矩阵无穷范数条件数 $\kappa(A)_\infty$. 显然

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

可以直接计算, 主要的困难在于计算 $\|A^{-1}\|_{\infty}$.

下面介绍的算法是由 Cline 等^[24] 给出的, 它在列主元 LU 分解 $PA = LU$ 的基础上, 再增加 $O(n^2)$ 的工作量, 就可以给出条件数的量级上的估计, 关键的是这个过程中不要求 A^{-1} .

设 y 为方程组 $Ay = d$ 的解, 则

$$\|A^{-1}\|_{\infty} \geq \frac{\|y\|_{\infty}}{\|d\|_{\infty}}.$$

上式表明, 线性方程组的解和右端项在无穷范数意义下的比是 $\|A^{-1}\|_{\infty}$ 的一个下界. 如果能选取向量 d 使得 $\|y\|_{\infty}/\|d\|_{\infty}$ 尽量大, 那么这个比值将会是 $\|A^{-1}\|_{\infty}$ 的一个很好的估计. 下面讨论如何选取 d . 设 A 的奇异值分解为

$$A = U\Sigma V^T,$$

式中: U 和 V 为正交阵; $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_n)$ 为由 A 的奇异值构成的对角矩阵.

不失一般性, 设 $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n > 0$. 设 d 为方程组 $A^T d = b$ 的解, $b \in \mathbb{R}^n$. 令 $V^T b = (\alpha_1, \alpha_2, \dots, \alpha_n)^T$, 则有

$$\|d\|_2^2 = \|A^{-T}b\|_2^2 = \|U\Sigma^{-1}V^Tb\|_2^2 = \|\Sigma^{-1}V^Tb\|_2^2 = \sum_{i=1}^n (\sigma_i^{-1}\alpha_i)^2, \quad (6.38)$$

及

$$\|y\|_2^2 = \|A^{-1}d\|_2^2 = \|A^{-1}(A^{-T}b)\|_2^2 = \|V\Sigma^{-2}V^Tb\|_2^2 = \|\Sigma^{-2}V^Tb\|_2^2 = \sum_{i=1}^n (\sigma_i^{-2}\alpha_i)^2, \quad (6.39)$$

注意到

$$\sum_{i=1}^n (\sigma_i^{-2}\alpha_i)^2 \leq \frac{1}{\sigma_n^2} \sum_{i=1}^n (\sigma_i^{-1}\alpha_i)^2,$$

从而有

$$\frac{1}{\sigma_n} \geq \frac{\|y\|_2}{\|d\|_2}.$$

从上式可知, $1/\sigma_n$ 是 $\|y\|_2/\|d\|_2$ 的上界, 且当 α_n 越大时, $\|y\|_2/\|d\|_2$ 越接近 $1/\sigma_n (= \|A^{-1}\|_2)$. 利用范数的等价性, $\|y\|_{\infty}/\|d\|_{\infty}$ 也会越大. 假设 A 有列主元 LU 分解 $PA = LU$. 令 $\tilde{d} = Pd$, 那么有

$$A^T d = U^T L^T \tilde{d} = b,$$

而方程组 $Ay = d$ 等价于

$$PAy = LUy = Pd = \tilde{d},$$

注意到 $\|Pd\|_{\infty} = \|d\|_{\infty}$, 则计算 $\|y\|_{\infty}/\|d\|_{\infty}$ 可以通过下面几步求得:

- (1) 求解 $U^T x = b$ 和 $L^T \tilde{d} = x$.
- (2) 求解 $Lz = \tilde{d}$ 和 $Uy = z$.

(3) 计算 $\|y\|_\infty/\|\tilde{d}\|_\infty$.

现考虑向量 b 的选取. 由式 (6.38) 可知, 若 α_n 越大, $\|d\|_2$ 也越大, 所以可以考虑选取向量 b 使得 x 的无穷范数尽可能大. 此外, 由于 b 乘上一个非零常数对 $\|y\|_\infty/\|d\|_\infty$ 的值没有影响, 不妨要求 b 的分量为 1 或 -1. 求解 $U^T x = b$ 的过程可通过下面算法给出:

```

 $x_1 = b_1/u_{11}; \beta_1 = 0;$ 
for  $k = 2:n$ 
     $\beta_k = \sum_{i=1}^{k-1} u_{ik}x_i; x_k = (b_k - \beta_k)/u_{kk};$ 
end

```

为了使得 $\|x\|_\infty$ 尽可能大, 取 $b_k = -\text{sign}(\beta_k)$, 则

$$|x_k| = \frac{|\text{sign}(\beta_k) + \beta_k|}{|u_{kk}|} = \frac{1 + |\beta_k|}{|u_{kk}|},$$

所以当 $|\beta_k|$ 越大时, $|x_k|$ 也会越大.

矩阵无穷范数条件数估算的 MATLAB 程序如下:

```

function [kappa]=cond_inf(A)
%输入: 矩阵A
%输出: 矩阵A无穷范数条件数的估计值
n=size(A,1); [L,U]=lu(A); %列主元LU分解PA=LU
%求解U'x=b
beta(1)=0; x=zeros(n,1); b(1)=1; x(1)=b(1)/U(1,1);
for k=2:n
    beta(k)=U(1:k-1,k)'*x(1:k-1);
    a=beta(k); b(k)=-sign(a);
    x(k)=(b(k)-beta(k))/U(k,k);
end
dw=L'\x; z=L\dw; y=U\z;
h2=norm(y,'inf')/norm(dw,'inf');
h1=norm(A,'inf'); kappa=h1*h2; %矩阵A的无穷范数条件数

```

例 6.14 利用程序 cond_inf.m 估算下面矩阵 A 的无穷范数条件数 $\kappa(A)_\infty$:

$$A = \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{bmatrix}_{n \times n}.$$

解 取 $n = 128$, 在 MATLAB 命令窗口依次输入:

```
>> n=128;e=ones(n-1,1);A=2*eye(n)+diag(e,-1)+diag(e,1);
>> [kappa]=cond_inf(A)
```

即可得到数值结果 (略).

6.5.3 舍入误差对解的影响

用直接法解线性方程组 $Ax = b$, 由于测量误差或舍入误差的存在, 导致 A 和 b 有误差 δA 和 δb . δA 和 δb 同计算机运算和精度有关, 计算精度越高, $\|\delta A\|$ 和 $\|\delta b\|$ 必然越小. 下面分析 $\|\delta A\|$ 和 $\|\delta b\|$ 对解的扰动的影响.

定义 6.5 如果矩阵 A 或 b 的微小变化, 引起方程组 $Ax = b$ 的解的巨大变化, 则称方程组是病态的, 称矩阵 A 为病态矩阵. 否则, 称方程组是良态的, 称矩阵 A 为良态矩阵.

当 A 和 b 都有微小的变化时, 解的变化可以用下面定理描述.

定理 6.6 设 $A \in \mathbb{R}^{n \times n}$ 为非奇异矩阵, $\delta A \in \mathbb{R}^{n \times n}$ 满足 $\|\delta A\| \cdot \|A^{-1}\| < 1$. 若 x 和 δx 满足

$$Ax = b, \quad (A + \delta A)(x + \delta x) = b + \delta b, \quad (6.40)$$

则

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\kappa(A)}{1 - \kappa(A)\varepsilon_r(A)} (\varepsilon_r(A) + \varepsilon_r(b)), \quad (6.41)$$

式中: $\varepsilon_r(A) = \frac{\|\delta A\|}{\|A\|}$, $\varepsilon_r(b) = \frac{\|\delta b\|}{\|b\|}$, 这里的矩阵范数是由向量范数诱导出来的算子范数.

证明 由于 $A + \delta A = (I + \delta A A^{-1})A$ 和 $\|\delta A A^{-1}\| \leq \|\delta A\| \cdot \|A^{-1}\| < 1$, 则 $A + \delta A$ 可逆且

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

可知

$$\|(A + \delta A)^{-1}\| = \|A^{-1}(I + \delta A A^{-1})^{-1}\| \leq \frac{\|A^{-1}\|}{1 - \|\delta A A^{-1}\|}.$$

此外, 由式 (6.40), 得

$$\begin{aligned} \delta x &= (A + \delta A)^{-1}(b + \delta b) - x \\ &= (A + \delta A)^{-1}b + (A + \delta A)^{-1}\delta b - A^{-1}b \\ &= ((A + \delta A)^{-1} - A^{-1})b + (A + \delta A)^{-1}\delta b \\ &= -(A + \delta A)^{-1}\delta A A^{-1}b + (A + \delta A)^{-1}\delta b \\ &= -(A + \delta A)^{-1}\delta A x + (A + \delta A)^{-1}\delta b. \end{aligned}$$

从而

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

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

两边除以 $\|x\|$, 注意到 $\|b\| = \|Ax\| \leq \|A\| \cdot \|x\|$, 可得

$$\begin{aligned} \frac{\|\delta x\|}{\|x\|} &\leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \cdot \|\delta A\|} \left(\|\delta A\| + \frac{\|\delta b\|}{\|x\|} \right) \\ &= \frac{\|A^{-1}\| \cdot \|A\|}{1 - \|A^{-1}\| \cdot \|A\| \cdot \frac{\|\delta A\|}{\|A\|}} \left(\frac{\|\delta A\|}{\|A\|} + \frac{\|\delta b\|}{\|A\| \cdot \|x\|} \right) \\ &\leq \frac{\|A^{-1}\| \cdot \|A\|}{1 - \|A^{-1}\| \cdot \|A\| \cdot \frac{\|\delta A\|}{\|A\|}} \left(\frac{\|\delta A\|}{\|A\|} + \frac{\|\delta b\|}{\|b\|} \right), \end{aligned}$$

即得到式 (6.41). 证毕. \square

根据定理 6.6, 当 $\kappa(A)$ 越大时, A 和 b 扰动之后得到的解的相对误差也越大, 相应的方程组也越病态.

对于病态的方程组, 为了得到较精确的近似解, 可以采用如下措施来减少舍入误差的影响: ① 采用高精度计算; ② 采用数值稳定性较好的算法; ③ 采用迭代改善计算解的办法.

设 x 是方程组 $Ax = b$ 的一个近似解, 如果它没有达到精度要求, 怎样产生一个更好的近似解呢? 记 $r = b - Ax$, 由于 x 为近似解, 所以 $r \neq 0$. 考虑方程组 $Az = r$, 如果 z 是 $Az = r$ 的精确解, 则 $A(x + z) = Ax + r = b$, 即 $x + z$ 为 $Ax = b$ 的一个精确解. 实际中 z 可能还是 $Az = r$ 的近似解. 当 $x + z$ 还是不够精确时, 把 $x + z$ 看作上述的 x , 重复上述过程. 设矩阵 A 有列主元 LU 分解 $PA = LU$, 下面给出对近似解的迭代改进:

- (1) $r := b - Ax$; (2) 解方程组 $Ly = Pr$, 得到 y ;
- (3) 解方程组 $Uz = y$, 得到 z ; (4) 令 $x := x + z$;
- (5) 若 x 达到精度的要求, 停算; 否则, 转 (1).

习题 6

6.1 用列主元 Gauss 消去法解下面的方程组

$$(1) \begin{cases} -3x_1 + 2x_2 + 6x_3 = 4, \\ 10x_1 - 7x_2 = 7, \\ 5x_1 - x_2 + 5x_3 = 6; \end{cases} \quad (2) \begin{cases} x_1 + 2x_2 + 3x_3 = 6, \\ 5x_1 - 6x_2 + 9x_3 = 8, \\ 3x_1 - 2x_2 + x_3 = 2. \end{cases}$$

6.2 顺序 Gauss 消去法可行的条件是 $a_{11}^{(1)}, a_{22}^{(2)}, \dots, a_{n-1,n-1}^{(n-1)}$ 都不为零. 试证明顺序 Gauss 消去法可行的充要条件是 A 的顺序主子式 $D_k \neq 0, 1 \leq k \leq n$.

6.3 设 $A = (a_{ij})$, $a_{11} \neq 0$, 经一步 Gauss 消元, 得

$$A^{(2)} = \begin{bmatrix} a_{11} & \alpha_1^T \\ O & A_2 \end{bmatrix}, \quad \text{其中 } A_2 = \begin{bmatrix} a_{22}^{(2)} & \cdots & a_{2n}^{(2)} \\ \vdots & \ddots & \vdots \\ a_{n2}^{(2)} & \cdots & a_{nn}^{(2)} \end{bmatrix}.$$

试证明:

- (1) 若 A 对称, 则 A_2 也对称;
- (2) 若 A 对称正定, 则 A_2 也对称正定.

6.4 对于 n 阶矩阵 $A = (a_{ij})$, 若

$$|a_{ii}| > \sum_{j=1, j \neq i}^n |a_{ij}|, \quad i = 1, 2, \dots, n,$$

则称 A 是严格对角占优矩阵. 证明: 若 A 是严格对角占优矩阵, 则经一步顺序 Gauss 消元过程后, 得到的 $A^{(1)}$ 仍为严格对角占优矩阵.

6.5 已知方程组 $Ax = f$, 其中

$$A = \begin{bmatrix} 2 & -1 & b \\ -1 & 2 & a \\ b & -1 & 2 \end{bmatrix}, \quad f = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}.$$

- (1) 试问参数满足什么条件式时, 可选用 Cholesky 分解法求解该方程组?
- (2) 取 $b = 0$, $a = 1$, 试用追赶法求解该方程组.

6.6 试证明:

- (1) 正定矩阵必存在 LU 分解;
- (2) 如果对称矩阵的各阶顺序主子式不等于零, 则必存在 LU 分解.

6.7 证明: 非奇异矩阵 A 不一定有 LU 分解.

6.8 证明: 非奇异矩阵 $A \in \mathbf{R}^{n \times n}$ 有唯一 LDU 分解的充要条件是 A 的顺序主子式 D_1, D_2, \dots, D_{n-1} 都是非零的.

6.9 已知方程组

$$\begin{bmatrix} 1 & 0 & -1 \\ 2 & 2 & 1 \\ 0 & 2 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1/2 \\ 1/3 \\ -2/3 \end{bmatrix}$$

的解为 $x = (1/2, -1/3, 0)^T$. 如果右端有微小扰动 $\|\delta b\|_\infty = 0.5 \times 10^{-6}$, 估计由此引起的解的相对误差.

6.10 方程组 $Ax = b$, 其中 A 为 $m \times n$ 阶对称且非奇异矩阵. 设 A 有误差 δA , 则原方程组变化为 $(A + \delta A)(x + \delta x) = b$, 其中 δx 为解的误差向量. 证明:

$$\frac{\|\delta x\|_2}{\|x + \delta x\|_2} \leq \left| \frac{\lambda_1}{\lambda_n} \right| \frac{\|\delta A\|_2}{\|A\|_2},$$

式中: λ_1 和 λ_n 分别为 A 的按模最大和最小的特征值.

第 7 章 矩阵特征值问题的数值方法

许多工程实际问题的求解,如振动问题、稳定性问题等,最终都归结为求某些矩阵的特征值和特征向量的问题. n 阶方阵 $A \in \mathbb{C}^{n \times n}$ 的特征值与特征向量,是满足如下两个方程的数 $\lambda \in \mathbb{C}$ 和非零向量 $x \in \mathbb{C}^n$:

$$p(\lambda) = \det(A - \lambda I) = 0, \quad (7.1)$$

$$Ax = \lambda x \text{ 或 } (A - \lambda I)x = 0. \quad (7.2)$$

式 (7.1) 称为矩阵 A 的特征方程, I 是 n 阶单位阵, $\det(A - \lambda I)$ 表示方阵 $A - \lambda I$ 的行列式,它是 λ 的 n 次代数多项式,当 n 较大时其零点难以准确求解.因此,从数值计算的观点来看,用特征多项式来求矩阵特征值的方法并不可取,必须建立有效的数值方法.

在实际应用中,求矩阵的特征值和特征向量通常采用迭代法.其基本思想是,将特征值和特征向量作为一个无限序列的极限来求得.舍入误差对这类方法的影响很小,但通常计算量较大.

根据具体问题的需要,有些实际问题只要计算模最大的特征值.当然,更多的问题则要求计算全部特征值和特征向量.本章介绍几种目前在计算机上比较常用的矩阵特征值问题的数值方法.

7.1 矩阵的特征值估计和隔离

在数值计算和其他学科中,往往需要估计一个矩阵的特征值在复平面上的位置.例如,在研究一个迭代算法时,需要判断迭代矩阵(或迭代函数的 Jacobi 矩阵)的特征值是否全部落在单位圆内.又如,分析一个非线性系统是否稳定时,需要知道有关矩阵的特征值是否全部落在复平面的左半部.本节扼要介绍一些这方面的结果.

首先讨论 Hermite 矩阵的特征值表示问题.由于 Hermite 矩阵 $A \in \mathbb{C}^{n \times n}$ 的特征值均为实数,故可约定其 n 个特征值的排列次序为:

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n. \quad (7.3)$$

定义 7.1 设 $A \in \mathbb{C}^{n \times n}$ 为 Hermite 矩阵,对任意的非零向量 $x \in \mathbb{C}^n$,称

$$R(x) = \frac{(Ax, x)}{(x, x)}$$

为 x 的 Rayleigh 商.

下面的定理说明了 Hermite 矩阵的特征值可以用 Rayleigh 商的极大值和极小值来描述.

定理 7.1 设矩阵 $A^H = A \in \mathbb{C}^{n \times n}$, 则

$$\max_{0 \neq x \in \mathbb{C}^n} \frac{(Ax, x)}{(x, x)} = \lambda_{\max}(A), \quad \min_{0 \neq x \in \mathbb{C}^n} \frac{(Ax, x)}{(x, x)} = \lambda_{\min}(A). \quad (7.4)$$

证明 设 A 的特征值为 $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$, 则存在酉矩阵 Q 使得

$$A = Q \operatorname{diag}(\lambda_1, \lambda_2, \cdots, \lambda_n) Q^H,$$

故

$$\begin{aligned} \frac{(Ax, x)}{(x, x)} &= \frac{(Q \operatorname{diag}(\lambda_1, \lambda_2, \cdots, \lambda_n) Q^H x, x)}{(Q^H x, Q^H x)} \\ &= \frac{(\operatorname{diag}(\lambda_1, \lambda_2, \cdots, \lambda_n) Q^H x, Q^H x)}{(Q^H x, Q^H x)}. \end{aligned}$$

令

$$y = Q^H x = (y_1, y_2, \cdots, y_n)^T,$$

则

$$\begin{aligned} \frac{(Ax, x)}{(x, x)} &= \frac{(\operatorname{diag}(\lambda_1, \lambda_2, \cdots, \lambda_n) y, y)}{(y, y)} \\ &= \frac{\lambda_1 y_1^2 + \lambda_2 y_2^2 + \cdots + \lambda_n y_n^2}{y_1^2 + y_2^2 + \cdots + y_n^2}. \end{aligned}$$

由于

$$\lambda_n(y_1^2 + y_2^2 + \cdots + y_n^2) \leq \lambda_1 y_1^2 + \lambda_2 y_2^2 + \cdots + \lambda_n y_n^2 \leq \lambda_1(y_1^2 + y_2^2 + \cdots + y_n^2),$$

故

$$\lambda_n \leq \frac{(Ax, x)}{(x, x)} \leq \lambda_1, \quad \forall 0 \neq x \in \mathbb{C}^n.$$

取 x_1 和 x_n 分别为对应于 λ_1 和 λ_n 的特征向量, 则

$$\frac{(Ax_1, x_1)}{(x_1, x_1)} = \lambda_1, \quad \frac{(Ax_n, x_n)}{(x_n, x_n)} = \lambda_n,$$

因此结论成立. 证毕. □

定理 7.2 $A^H = A \in \mathbb{C}^{n \times n}$ 的特征值如式 (7.3), 则对 $1 \leq k \leq n$, 有

$$\lambda_k = \max_{V_k} \min_{0 \neq x \in V_k} \frac{(Ax, x)}{(x, x)}, \quad (7.5)$$

式中: V_k 为 \mathbb{C}^n 的任意一个 k 维子空间.

证明 设 A 的对应于特征值 $\lambda_1, \lambda_2, \cdots, \lambda_n$ 的特征向量依次为 p_1, p_2, \cdots, p_n , 且标准正交, 构造子空间 $W_k = \operatorname{span}\{p_k, p_{k+1}, \cdots, p_n\}$, 那么 $\dim(W_k) = n - k + 1$. 由于 $V_k + W_k \subset \mathbb{C}^n$, 利用子空间的维数公式求得

$$n \geq \dim(V_k + W_k) = \dim(V_k) + \dim(W_k) - \dim(V_k \cap W_k)$$

$$= n + 1 - \dim(\mathcal{V}_k \cap \mathcal{W}_k),$$

即 $\dim(\mathcal{V}_k \cap \mathcal{W}_k) \geq 1$. 于是存在 $\mathbf{x}_0 \in \mathcal{V}_k \cap \mathcal{W}_k \subset \mathcal{W}_k$ 满足 $\|\mathbf{x}_0\|_2 = 1$, 且有

$$\mathbf{x}_0 = c_k \mathbf{p}_k + c_{k+1} \mathbf{p}_{k+1} + \cdots + c_n \mathbf{p}_n, \quad (c_k^2 + c_{k+1}^2 + \cdots + c_n^2 = 1),$$

$$(\mathbf{A}\mathbf{x}_0, \mathbf{x}_0) = \lambda_k c_k^2 + \lambda_{k+1} c_{k+1}^2 + \cdots + \lambda_n c_n^2 \leq \lambda_k.$$

因此

$$\min_{0 \neq \mathbf{x} \in \mathcal{V}_k} \frac{(\mathbf{A}\mathbf{x}, \mathbf{x})}{(\mathbf{x}, \mathbf{x})} \leq \frac{(\mathbf{A}\mathbf{x}_0, \mathbf{x}_0)}{(\mathbf{x}_0, \mathbf{x}_0)} \leq \lambda_k.$$

由 \mathcal{V}_k 的任意性, 得

$$\max_{\mathcal{V}_k} \min_{0 \neq \mathbf{x} \in \mathcal{V}_k} \frac{(\mathbf{A}\mathbf{x}, \mathbf{x})}{(\mathbf{x}, \mathbf{x})} \leq \lambda_k.$$

另外, 取 k 维子空间 $\mathcal{V}_k = \text{span}\{\mathbf{p}_1, \mathbf{p}_2, \cdots, \mathbf{p}_k\}$, 当 $\mathbf{x} \in \mathcal{V}_k$ 满足 $\|\mathbf{x}\|_2 = 1$ 时, 有

$$\mathbf{x} = c_1 \mathbf{p}_1 + c_2 \mathbf{p}_2 + \cdots + c_k \mathbf{p}_k, \quad (c_1^2 + c_2^2 + \cdots + c_k^2 = 1),$$

$$(\mathbf{A}\mathbf{x}, \mathbf{x}) = \lambda_1 c_1^2 + \lambda_2 c_2^2 + \cdots + \lambda_k c_k^2 \geq \lambda_k.$$

由 $\mathbf{x} \in \mathcal{V}_k$ 的任意性, 得

$$\min_{0 \neq \mathbf{x} \in \mathcal{V}_k} \frac{(\mathbf{A}\mathbf{x}, \mathbf{x})}{(\mathbf{x}, \mathbf{x})} \geq \lambda_k.$$

从而有

$$\max_{\mathcal{V}_k} \min_{0 \neq \mathbf{x} \in \mathcal{V}_k} \frac{(\mathbf{A}\mathbf{x}, \mathbf{x})}{(\mathbf{x}, \mathbf{x})} \geq \lambda_k.$$

综上所述, 式 (7.5) 成立. 证毕. □

定理 7.3 设 n 阶 Hermite 矩阵 \mathbf{A} 和 $\mathbf{B} = \mathbf{A} + \mathbf{E}$ 的特征值依次为

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n, \quad \mu_1 \geq \mu_2 \geq \cdots \geq \mu_n.$$

则有

$$\varepsilon_n \leq \mu_k - \lambda_k \leq \varepsilon_1, \tag{7.6}$$

式中: ε_1 和 ε_n 分别为 Hermite 矩阵 \mathbf{E} 的最大特征值和最小特征值.

证明 下面约定向量 \mathbf{x} 满足 $\|\mathbf{x}\|_2 = 1$. 在式 (7.5) 中固定 \mathcal{V}_k , 利用式 (7.4), 得

$$\begin{aligned} \mu_k &\geq \min_{\mathbf{x} \in \mathcal{V}_k} (\mathbf{B}\mathbf{x}, \mathbf{x}) \geq \min_{\mathbf{x} \in \mathcal{V}_k} (\mathbf{A}\mathbf{x}, \mathbf{x}) + \min_{\mathbf{x} \in \mathcal{V}_k} (\mathbf{E}\mathbf{x}, \mathbf{x}) \\ &\geq \min_{\mathbf{x} \in \mathcal{V}_k} (\mathbf{A}\mathbf{x}, \mathbf{x}) + \min_{\mathbf{x} \in \mathbb{C}^n} (\mathbf{E}\mathbf{x}, \mathbf{x}) = \min_{\mathbf{x} \in \mathcal{V}_k} (\mathbf{A}\mathbf{x}, \mathbf{x}) + \varepsilon_n, \\ \lambda_k &\geq \min_{\mathbf{x} \in \mathcal{V}_k} (\mathbf{A}\mathbf{x}, \mathbf{x}) \geq \min_{\mathbf{x} \in \mathcal{V}_k} (\mathbf{B}\mathbf{x}, \mathbf{x}) + \min_{\mathbf{x} \in \mathcal{V}_k} (-\mathbf{E}\mathbf{x}, \mathbf{x}) \\ &\geq \min_{\mathbf{x} \in \mathcal{V}_k} (\mathbf{B}\mathbf{x}, \mathbf{x}) + \min_{\mathbf{x} \in \mathbb{C}^n} (-\mathbf{E}\mathbf{x}, \mathbf{x}) = \min_{\mathbf{x} \in \mathcal{V}_k} (\mathbf{B}\mathbf{x}, \mathbf{x}) + (-\varepsilon_1). \end{aligned}$$

由于 $\mathcal{V}_k \subset \mathbb{R}^n$ 任意, 所以

$$\begin{aligned}\mu_k &\geq \max_{\mathcal{V}_k} \min_{x \in \mathcal{V}_k} (Ax, x) + \varepsilon_n = \lambda_k + \varepsilon_n, \\ \lambda_k &\geq \max_{\mathcal{V}_k} \min_{x \in \mathcal{V}_k} (Bx, x) - \varepsilon_1 = \mu_k - \varepsilon_1,\end{aligned}$$

即式 (7.6) 成立. 证毕. \square

注 7.1 定理 7.3 可用来讨论实对称矩阵的摄动问题. 可将 E 看作由于种种原因在矩阵 A 的元素中产生的误差构成的“误差矩阵”, 而 B 就是实际计算时的矩阵. 设对 E 的每个元素 e_{ij} 都有 $|e_{ij}| \leq \varepsilon$, 其中 $\varepsilon > 0$ 是某个常数, 则 B 的特征值 $\lambda(B) = \lambda(A) + \Delta\lambda$ 与 A 的特征值 $\lambda(A)$ 之间满足 (7.6) 式, 即 $\varepsilon_1 \geq \Delta\lambda \geq \varepsilon_n$. 因此, 有

$$|\Delta\lambda| \leq \rho(E) \leq \|E\|_F \leq n\varepsilon.$$

这表明, 当 A 的诸元素有一个小扰动时, 在 A 的特征值中造成的扰动与其同阶.

定理 7.4 (分隔定理) 设 A 是 n 阶实对称矩阵, $B = Q^T A Q$, 其中 $Q \in \mathbb{R}^{n \times (n-1)}$ 满足 $Q^T Q = I_{n-1}$. 再设 A 和 B 的特征值分别为

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n \quad \text{和} \quad \mu_1 \geq \mu_2 \geq \cdots \geq \mu_{n-1},$$

则有

$$\lambda_1 \geq \mu_1 \geq \lambda_2 \geq \mu_2 \geq \cdots \geq \mu_{n-1} \geq \lambda_n. \quad (7.7)$$

特别地, 在定理 7.4 中选取 $Q = [e_1, \cdots, e_{i-1}, e_{i+1}, \cdots, e_n]$, 即得如下结论.

推论 7.1 设 B 是 n 阶实对称矩阵 A 的一个 $n-1$ 阶主子阵, 并假定 A 和 B 的特征值分别为

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n \quad \text{和} \quad \mu_1 \geq \mu_2 \geq \cdots \geq \mu_{n-1},$$

则有

$$\lambda_1 \geq \mu_1 \geq \lambda_2 \geq \mu_2 \geq \cdots \geq \mu_{n-1} \geq \lambda_n.$$

反复应用推论 7.1, 便有以下推论.

推论 7.2 设 A 是一个 n 阶实对称矩阵, B 是 A 的一个 k 阶主子阵 ($1 \leq k \leq n-1$), 并假定 A 和 B 的特征值分别为

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n \quad \text{和} \quad \mu_1 \geq \mu_2 \geq \cdots \geq \mu_k,$$

则有

$$\lambda_i \geq \mu_i \geq \lambda_{n-k+i}, \quad i = 1, 2, \cdots, k. \quad (7.8)$$

为了细致描述 n 阶矩阵的特征值在复平面上的分布范围, 下面引进 Gerschgorin 圆盘 (简称盖尔圆或盖氏圆).

定义 7.2 设 $A = (a_{ij}) \in \mathbb{C}^{n \times n}$, 令 $R_i = \sum_{j=1, j \neq i}^n |a_{ij}|$, 则称

$$G_i = \{z | z \in \mathbb{C} : |z - a_{ii}| \leq R_i\}, \quad i = 1, 2, \dots, n \quad (7.9)$$

为 A 的第 i 个盖氏圆.

定理 7.5 设 λ 是 $A \in \mathbb{C}^{n \times n}$ 的任一特征值, 则 $\lambda \in \bigcup_{i=1}^n G_i$, 即 A 的任一特征值都落在它的 n 个盖氏圆盘的并集内.

证明 设 A 的对应于特征值 λ 的特征向量为 $x = (x_1, x_2, \dots, x_n)^T$. 选取 i_0 使得 $|x_{i_0}| = \max_{1 \leq i \leq n} |x_i|$, 则由 $Ax = \lambda x$ 可得

$$\begin{aligned} \sum_{j=1}^n a_{i_0 j} x_j &= \lambda x_{i_0} \implies (\lambda - a_{i_0 i_0}) x_{i_0} = \sum_{j=1, j \neq i_0}^n a_{i_0 j} x_j \\ \implies |\lambda - a_{i_0 i_0}| &= \left| \sum_{j=1, j \neq i_0}^n a_{i_0 j} \frac{x_j}{x_{i_0}} \right| \leq R_{i_0}, \end{aligned}$$

即 $\lambda \in G_{i_0} \subset \bigcup_{i=1}^n G_i$. 证毕. \square

定理 7.5 用一组圆盘覆盖矩阵的特征值分布区域, 下面介绍用另外一组几何图形覆盖矩阵的特征值分布区域的定理, 后者可以看作是前者的推广.

定理 7.6 设 λ 是 $A \in \mathbb{C}^{n \times n}$ ($n > 1$) 的任一特征值, 则 λ 位于某个

$$\Omega_{ij} = \{z | z \in \mathbb{C}, |z - a_{ii}| |z - a_{jj}| \leq R_i R_j, i \neq j; i, j = 1, 2, \dots, n\}$$

之中, 称 Ω_{ij} ($i \neq j$) 为 A 的 Cassini (卡西尼) 卵形.

证明 设 A 的对应于特征值 λ 的特征向量为 $x = (x_1, x_2, \dots, x_n)^T$. 选取 $i_0 \neq j_0$ 满足 $|x_{i_0}| \geq |x_{j_0}| \geq |x_k|$ ($k \neq i_0, j_0$), 下证 $\lambda \in \Omega_{i_0 j_0}$.

(1) 如果 $x_{j_0} = 0$, 则 $x_{i_0} \neq 0$, $x_k = 0$ ($k \neq i_0$). 由 $Ax = \lambda x$ 可得

$$\lambda x_{i_0} = \sum_{k=1}^n a_{i_0 k} x_k = a_{i_0 i_0} x_{i_0} \Rightarrow \lambda = a_{i_0 i_0}.$$

故

$$|\lambda - a_{i_0 i_0}| |\lambda - a_{j_0 j_0}| = 0 \leq R_{i_0} R_{j_0}.$$

(2) 如果 $x_{j_0} \neq 0$, 则 $|x_{i_0}| \geq |x_{j_0}| > 0$, 再由 $Ax = \lambda x$ 可得

$$(\lambda - a_{ii}) x_i = \sum_{k \neq i} a_{ik} x_k, \quad (i = 1, 2, \dots, n).$$

取 $i = i_0$ 时, 可得

$$|\lambda - a_{i_0 i_0}| |x_{i_0}| \leq \sum_{k \neq i_0} |a_{i_0 k}| |x_k| \leq |x_{j_0}| R_{i_0}.$$

取 $i = j_0$ 时, 可得

$$|\lambda - a_{j_0 j_0}| |x_{j_0}| \leq \sum_{k \neq j_0} |a_{j_0 k}| |x_k| \leq |x_{i_0}| R_{j_0}.$$

因此 $|\lambda - a_{i_0 i_0}| |\lambda - a_{j_0 j_0}| \leq R_{i_0} R_{j_0}$. 综述 (1) 和 (2) 即得 $\lambda \in \Omega_{i_0 j_0}$. 证毕. \square

推论 7.3 设 $A = (a_{ij}) \in \mathbb{C}^{n \times n}$ ($n > 1$) 满足 $|a_{ii}| |a_{jj}| > R_i R_j$ ($i \neq j$), 则 $\det(A) \neq 0$.

证明 设 λ 是 A 的任一特征值, 那么必有 Ω_{ij} 使得 $\lambda \in \Omega_{ij}$, 即

$$|\lambda - a_{ii}| |\lambda - a_{jj}| \leq R_i R_j.$$

如果 $\lambda = 0$, 则有 $|a_{ii}| |a_{jj}| \leq R_i R_j$, 这与题设矛盾, 故 $\lambda \neq 0$. 从而 $\det(A) \neq 0$. 证毕. \square

7.2 幂法和反幂法

7.2.1 幂法

幂法是通过求矩阵的特征向量来求出特征值的一种迭代法. 它主要用来求按模最大的特征值和相应的特征向量的. 其优点是算法简单, 容易计算机实现, 缺点是收敛速度慢, 其有效性依赖于矩阵特征值的分布情况.

适于使用幂法的常见情形是: A 的特征值可按模的大小排列为 $|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|$, 且其对应特征向量 $\xi_1, \xi_2, \dots, \xi_n$ 线性无关. 此时, 任意非零向量 $x^{(0)}$ 均可用 $\xi_1, \xi_2, \dots, \xi_n$ 线性表示, 即

$$x^{(0)} = \alpha_1 \xi_1 + \alpha_2 \xi_2 + \dots + \alpha_n \xi_n, \quad (7.10)$$

且 $\alpha_1, \alpha_2, \dots, \alpha_n$ 不全为零. 作向量序列 $x^{(k)} = A^k x^{(0)}$, 则

$$\begin{aligned} x^{(k)} &= A^k x^{(0)} = \alpha_1 A^k \xi_1 + \alpha_2 A^k \xi_2 + \dots + \alpha_n A^k \xi_n \\ &= \alpha_1 \lambda_1^k \xi_1 + \alpha_2 \lambda_2^k \xi_2 + \dots + \alpha_n \lambda_n^k \xi_n \\ &= \lambda_1^k \left[\alpha_1 \xi_1 + \alpha_2 \left(\frac{\lambda_2}{\lambda_1} \right)^k \xi_2 + \dots + \alpha_n \left(\frac{\lambda_n}{\lambda_1} \right)^k \xi_n \right]. \end{aligned}$$

由此可见, 若 $\alpha_1 \neq 0$, 则因 $k \rightarrow \infty$ 时, 有

$$\left(\frac{\lambda_i}{\lambda_1} \right)^k \rightarrow 0, \quad (i = 2, \dots, n),$$

故当 k 充分大时, 必有

$$x^{(k)} \approx \lambda_1^k \alpha_1 \xi_1,$$

即 $x^{(k)}$ 可以近似看成 λ_1 对应的特征向量; 而 $x^{(k)}$ 与 $x^{(k-1)}$ 分量之比为

$$\frac{x_i^{(k)}}{x_i^{(k-1)}} \approx \frac{\lambda_1^k \alpha_1 (\xi_1)_i}{\lambda_1^{k-1} \alpha_1 (\xi_1)_i} = \lambda_1.$$

于是利用向量序列 $\{x^{(k)}\}$ 既可求出按模最大的特征值 λ_1 , 又可求出对应的特征向量 ξ_1 .

在实际计算中, 考虑到当 $|\lambda_1| > 1$ 时, $\lambda_1^k \rightarrow \infty$; 当 $|\lambda_1| < 1$ 时, $\lambda_1^k \rightarrow 0$, 因而计算 $x^{(k)}$ 时可能会导致计算机“上溢”或“下溢”现象发生, 故采取每步将 $x^{(k)}$ 归一化处理的办法, 即将 $x^{(k)}$ 的各分量都除以模最大的分量, 使 $\|x^{(k)}\|_\infty = 1$. 于是, 求 A 按模最大的特征值 λ_1 和对应的特征向量 ξ_1 的算法, 可归纳为如下步骤.

算法 7.1 (幂法)

步 1, 输入矩阵 A , 初始向量 $v^{(0)}$, 误差限 ε , 最大迭代次数 N . 记 m_0 是 $v^{(0)}$ 按模最大的分量, $x^{(0)} = v^{(0)}/m_0$. 置 $k := 0$.

步 2, 计算 $v^{(k+1)} = Ax^{(k)}$. 记 m_{k+1} 是 $v^{(k+1)}$ 按模最大的分量, $x^{(k+1)} = v^{(k+1)}/m_{k+1}$.

步 3, 若 $|m_{k+1} - m_k| < \varepsilon$, 停算, 输出近似特征值 m_{k+1} 和近似特征向量 $x^{(k+1)}$; 否则, 转步 4.

步 4, 若 $k < N$, 置 $k := k + 1$, 转步 2; 否则输出计算失败信息, 停算.

算法 7.1 称为幂法. 幂法的算法结构简单, 容易编程实现. 其 MATLAB 程序如下:

```
function [lam,v,k]=mypower(A,x,tol,N)
%用幂法计算矩阵的模最大特征值和对应的特征向量
%输入:A为n阶方阵,x为初始向量,tol为控制精度,N为最大迭代次数
%输出:lam为按模最大的特征值,v为对应的特征向量,k为迭代次数
if nargin<4, N=1000; end
if nargin<3, tol=1e-6; end
m=0; k=0;
while(k<N)
    v=A*x;
    [m1,t]=max(abs(v));
    m1=v(t); x=v/m1;
    err=abs(m1-m);
    if err<tol, break; end
    m=m1; k=k+1;
end
lam=m1; v=x;
```

例 7.1 利用程序 mypower.m 求矩阵 A 按模最大的特征值 λ_1 和对应的特征向量 ξ_1 , 其中

$$A = \begin{bmatrix} 1 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 99 & -1 \\ & & & -1 & 100 \end{bmatrix}.$$

解 用幂法的 MATLAB 程序, 编写 MATLAB 脚本文件 ex71.m, 取容许误差为 $\varepsilon = 10^{-6}$, 在命令窗口运行之, 迭代 661 次得到模最大的特征值 $\lambda_1 = 100.7461$.

下面证明算法 7.1 的收敛性定理.

定理 7.7 设矩阵 A 的特征值可按模的大小排列为 $|\lambda_1| > |\lambda_2| \geq \cdots \geq |\lambda_n|$, 且其对应特征向量 $\xi_1, \xi_2, \cdots, \xi_n$ 线性无关. 序列 $\{x^{(k)}\}$ 由算法 7.1 产生, 则有

$$\lim_{k \rightarrow \infty} x^{(k)} = \frac{\xi_1}{\max\{\xi_1\}} := \xi_1^0, \quad \lim_{k \rightarrow \infty} m_k = \lambda_1, \quad (7.11)$$

式中: ξ_1^0 为将 ξ_1 归一化后得到的向量; $\max\{\xi_1\}$ 为向量 ξ_1 模最大的分量.

证明 由算法 7.1 的步 2 和步 3 知

$$x^{(k)} = \frac{v^{(k)}}{m_k} = \frac{Ax^{(k-1)}}{m_k} = \frac{A^2x^{(k-2)}}{m_k m_{k-1}} = \cdots = \frac{A^k x^{(0)}}{m_k m_{k-1} \cdots m_1}.$$

由于 $x^{(k)}$ 的最大分量为 1, 即 $\max\{x^{(k)}\} = 1$, 故

$$m_k m_{k-1} \cdots m_1 = \max\{A^k x^{(0)}\}.$$

从而

$$\begin{aligned} x^{(k)} &= \frac{A^k x^{(0)}}{\max\{A^k x^{(0)}\}} = \frac{\lambda_1^k \left[\alpha_1 \xi_1 + \sum_{i=2}^n \alpha_i \left(\frac{\lambda_i}{\lambda_1} \right)^k \xi_i \right]}{\max \left\{ \lambda_1^k \left[\alpha_1 \xi_1 + \sum_{i=2}^n \alpha_i \left(\frac{\lambda_i}{\lambda_1} \right)^k \xi_i \right] \right\}} \\ &= \frac{\alpha_1 \xi_1 + \sum_{i=2}^n \alpha_i \left(\frac{\lambda_i}{\lambda_1} \right)^k \xi_i}{\max \left\{ \alpha_1 \xi_1 + \sum_{i=2}^n \alpha_i \left(\frac{\lambda_i}{\lambda_1} \right)^k \xi_i \right\}}. \end{aligned}$$

可见

$$\lim_{k \rightarrow \infty} x^{(k)} = \frac{\alpha_1 \xi_1}{\max\{\alpha_1 \xi_1\}} = \frac{\xi_1}{\max\{\xi_1\}} = \xi_1^0.$$

又

$$\begin{aligned} v^{(k)} &= Ax^{(k-1)} = \frac{A^k x^{(0)}}{m_{k-1} \cdots m_1} = \frac{A^k x^{(0)}}{\max\{A^{k-1} x^{(0)}\}} \\ &= \frac{\lambda_1^k \left[\alpha_1 \xi_1 + \sum_{i=2}^n \alpha_i \left(\frac{\lambda_i}{\lambda_1} \right)^k \xi_i \right]}{\lambda_1^{k-1} \max \left\{ \alpha_1 \xi_1 + \sum_{i=2}^n \alpha_i \left(\frac{\lambda_i}{\lambda_1} \right)^{k-1} \xi_i \right\}}, \end{aligned}$$

注意到 m_k 是 $v^{(k)}$ 模最大的分量, 即有

$$m_k = \max\{v^{(k)}\} = \lambda_1 \frac{\max \left\{ \alpha_1 \xi_1 + \sum_{i=2}^n \alpha_i \left(\frac{\lambda_i}{\lambda_1} \right)^k \xi_i \right\}}{\max \left\{ \alpha_1 \xi_1 + \sum_{i=2}^n \alpha_i \left(\frac{\lambda_i}{\lambda_1} \right)^{k-1} \xi_i \right\}},$$

从而 $\lim_{k \rightarrow \infty} m_k = \lambda_1$ 成立. 证毕. □

7.2.2 幂法的加速技术

定理 7.8 在定理 7.7 的条件下, 幂法 (算法 7.1) 是线性收敛的.

证明 设 k 充分大时, $A^k x^{(0)}$ 的按模最大分量是它的第 j 个分量, 则

$$\begin{aligned} m_k - \lambda_1 &= \max\{v^{(k)}\} - \lambda_1 = \frac{\max\{A^k x^{(0)}\}}{\max\{A^{k-1} x^{(0)}\}} - \lambda_1 \\ &= \frac{[\beta_1 \lambda_1^k \xi_1 + \beta_2 \lambda_2^k \xi_2 + \cdots + \beta_n \lambda_n^k \xi_n]_j}{[\beta_1 \lambda_1^{k-1} \xi_1 + \beta_2 \lambda_2^{k-1} \xi_2 + \cdots + \beta_n \lambda_n^{k-1} \xi_n]_j} - \lambda_1 \\ &= \frac{[\beta_2 \lambda_2^{k-1} (\lambda_2 - \lambda_1) \xi_2 + \cdots + \beta_n \lambda_n^{k-1} (\lambda_n - \lambda_1) \xi_n]_j}{[\beta_1 \lambda_1^{k-1} \xi_1 + \beta_2 \lambda_2^{k-1} \xi_2 + \cdots + \beta_n \lambda_n^{k-1} \xi_n]_j}. \end{aligned}$$

于是有

$$\begin{aligned} m_k - \lambda_1 &= \left(\frac{\lambda_2}{\lambda_1}\right)^{k-1} \frac{\left[\beta_2 (\lambda_2 - \lambda_1) \xi_2 + \sum_{i=3}^n \beta_i \left(\frac{\lambda_i}{\lambda_2}\right)^{k-1} (\lambda_i - \lambda_1) \xi_i\right]_j}{\left[\beta_1 \xi_1 + \sum_{i=2}^n \beta_i \left(\frac{\lambda_i}{\lambda_1}\right)^{k-1} \xi_i\right]_j} \\ &= \left(\frac{\lambda_2}{\lambda_1}\right)^{k-1} M_k, \quad M_k \rightarrow M, \end{aligned}$$

式中: M 为常数. 所以, 当 $k \rightarrow \infty$ 时, 有

$$\frac{|m_{k+1} - \lambda_1|}{|m_k - \lambda_1|} = \left| \frac{M_{k+1} (\lambda_2 / \lambda_1)^k}{M_k (\lambda_2 / \lambda_1)^{k-1}} \right| \rightarrow \left| \frac{\lambda_2}{\lambda_1} \right|,$$

这就证明了幂法的线性收敛速度. 证毕. □

定理 7.8 表明, 幂法的收敛速度与比值 $|\lambda_2 / \lambda_1|$ 的大小有关, $|\lambda_2 / \lambda_1|$ 越小, 收敛速度越快, 当此比值接近于 1 时, 收敛速度是非常缓慢的. 因此, 可以对矩阵作一原点位移, 令

$$B = A - \alpha I,$$

式中: α 为参数. 选择此参数可使矩阵 B 的上述比值更小, 以加快幂法的收敛速度. 设矩阵 A 的特征值为 $\lambda_1, \lambda_2, \cdots, \lambda_n$, 对应的特征向量为 $\xi_1, \xi_2, \cdots, \xi_n$, 则矩阵 B 的特征值为 $\lambda_1 - \alpha, \lambda_2 - \alpha, \cdots, \lambda_n - \alpha$, B 的特征向量和 A 的特征向量相同. 假设原点位移后, B 的特征值 $\lambda_1 - \alpha$ 仍为模最大的特征值, 选择 α 的目的是使

$$\max_{2 \leq i \leq n} \frac{|\lambda_i - \alpha|}{|\lambda_1 - \alpha|} < \left| \frac{\lambda_2}{\lambda_1} \right|. \quad (7.12)$$

适当地选择 α 可使幂法的收敛速度得到加速. 此时 $m_k \rightarrow \lambda_1 - \alpha$, $m_k + \alpha \rightarrow \lambda_1$, 而 $x^{(k)}$ 仍然收敛于 A 的特征向量 ξ_1^0 . 这种加速收敛的方法称为原点位移法.

在实际计算中, 由于矩阵的特征值分布情况事先一般是不知道的, 参数 α 的选取存在困难, 因为 α 的选取要保证 $\lambda_1 - \alpha$ 仍然是矩阵 $B (= A - \alpha I)$ 模最大的特征值, 故原点位移法是很难实现的. 但是, 在反幂法中, 原点位移参数 α 是很容易选取的, 因此, 带原点位移的反幂法已成为改进特征值和特征向量精度的标准算法. 采用原点位移加速技术的幂法 MATLAB 程序如下:

```
function [lam,v,k]=mopower(A,x,alpha,tol,N)
%用原点位移幂法求矩阵的模最大特征值和对应的特征向量
%输入:A为n阶方阵,x为初始向量,tol为控制精度,
%      N最大迭代次数,alpha为原点位移参数
%输出:lam返回按模最大的特征值,v返回对应的特征向量,k返回迭代次数
if nargin<5, N=1000; end
if nargin<4, tol=1e-6; end
m=0; k=0;
A=A-alpha*eye(length(x));
while(k<N)
    v=A*x;
    [m1,t]=max(abs(v));
    m1=v(t); x=v/m1;
    err=abs(m1-m);
    if err<tol, break; end
    m=m1; k=k+1;
end
lam=m1+alpha;v=x;
```

例 7.2 利用原点位移幂法通用程序, 取 $\alpha = 50$, 求例 7.1 中的矩阵 A 按模最大的特征值 λ_1 和对应的特征向量 ξ_1 .

解 用原点位移幂法的 MATLAB 程序, 编写 M 文件 ex72.m, 在命令窗口运行该文件, 迭代 353 次得到模最大的特征值 $\lambda_1 = 100.7461$.

由计算结果可以看出, 在同样的精度控制下, 带原点位移加速的幂法只需要迭代 353 次, 而纯粹的幂法则需要迭代 661 次 (例 7.1).

7.2.3 反幂法

设 A 可逆, 则对 A 的逆阵 A^{-1} 施以幂法称为反幂法. 由于 $A\xi_i = \lambda_i\xi_i$ 时, 成立 $A^{-1}\xi_i = \lambda_i^{-1}\xi_i$. 因此, 若 $|\lambda_1| \geq |\lambda_2| \geq \cdots \geq |\lambda_{n-1}| > |\lambda_n|$, 则 λ_n^{-1} 是 A^{-1} 按模最大的特征值, 此时按反幂法, 必有

$$m_k \rightarrow \lambda_n^{-1}, \quad x^{(k)} \rightarrow \xi_n^0,$$

且其收敛率为 $|\lambda_n/\lambda_{n-1}|$. 任取初始向量 $\mathbf{x}^{(0)}$, 构造向量序列

$$\mathbf{x}^{(k+1)} = \mathbf{A}^{-1}\mathbf{x}^{(k)}, \quad k = 0, 1, 2, \dots \quad (7.13)$$

按幂法计算即可. 但用式 (7.13) 计算, 首先要求 \mathbf{A}^{-1} , 这比较麻烦而且是不经济的. 实际计算中, 通常用解方程组的办法, 即用

$$\mathbf{A}\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)}, \quad k = 0, 1, 2, \dots, \quad (7.14)$$

求 $\mathbf{x}^{(k+1)}$. 为防止计算机溢出, 实际计算时所用的公式为

$$\begin{cases} \mathbf{v}^{(k)} = \mathbf{x}^{(k)} / \max(\mathbf{x}^{(k)}), \\ \mathbf{A}\mathbf{x}^{(k+1)} = \mathbf{v}^{(k)}, \end{cases} \quad k = 0, 1, 2, \dots, \quad (7.15)$$

式中: $\max(\mathbf{x}^{(k)})$ 为 $\mathbf{x}^{(k)}$ 模最大的分量.

反幂法主要用于已知矩阵的近似特征值为 α 时, 求矩阵的特征向量并提高特征值的精度. 此时, 可以用原点位移法来加速迭代过程, 于是式 (7.15) 相应为

$$\begin{cases} \mathbf{v}^{(k)} = \mathbf{x}^{(k)} / \max(\mathbf{x}^{(k)}), \\ (\mathbf{A} - \alpha\mathbf{I})\mathbf{x}^{(k+1)} = \mathbf{v}^{(k)}, \end{cases} \quad k = 0, 1, 2, \dots. \quad (7.16)$$

反幂法的计算步骤如下.

算法 7.2 (反幂法)

步 1, 选取初值 $\mathbf{x}^{(0)}$, 近似值 α , 误差限 ε , 最大迭代次数 N . 记 m_0 为 $\mathbf{x}^{(0)}$ 中按模最大的分量, $\mathbf{v}^{(0)} = \mathbf{x}^{(0)} / m_0$. 置 $k := 0$.

步 2, 解方程组 $(\mathbf{A} - \alpha\mathbf{I})\mathbf{x}^{(k+1)} = \mathbf{v}^{(k)}$ 得 $\mathbf{x}^{(k+1)}$.

步 3, 记 m_{k+1} 为 $\mathbf{x}^{(k+1)}$ 中按模最大的分量, $\mathbf{v}^{(k+1)} = \mathbf{x}^{(k+1)} / m_{k+1}$.

步 4, 若 $|m_{k+1}^{-1} - m_k^{-1}| < \varepsilon$, 则置 $\lambda := m_{k+1}^{-1} + \alpha$, 输出 λ 和 $\mathbf{x}^{(k+1)}$, 停算; 否则, 转步 5.

步 5, 若 $k < N$, 置 $k := k + 1$, 转步 2, 否则输出计算失败信息, 停算.

注 7.2 (1) 算法 7.2 计算出与数 α 最接近的特征值及相应的特征向量. 若取 $\alpha = 0$, 则求出 \mathbf{A} 的按模最小的特征值.

(2) 通常首先用幂法求出 \mathbf{A} 的按模最大的近似特征值作为算法 7.2 中的 α 值, 再使用该算法对 α 和相应的特征向量进行精确化.

(3) 为节省计算量, 通常先用列主元 LU 分解将矩阵 $\mathbf{A} - \alpha\mathbf{I}$ 分解为下三角矩阵 \mathbf{L} 和上三角矩阵 \mathbf{U} , 这样在迭代过程中每一步就只需求解两个三角方程组即可.

反幂法的 MATLAB 程序如下:

```

function [lam,v,k]=mvpower(A,x,alpha,tol,N)
%用反幂法计算矩阵与alpha最接近的特征值和对应的特征向量
%输入:A为n阶方阵,x为初始向量,tol为精度,N为最大迭代数,alpha为某个常数
%输出:lam返回与alpha最接近的特征值,v返回对应的特征向量,k返回迭代次数
if nargin<5, N=500; end
if nargin<4, tol=1e-5; end
m=0.5; k=0;
A=A-alpha*eye(length(x));
[L,U,P]=lu(A);
while (k<N)
    [m1,t]=max(abs(x));
    m1=x(t); v=x/m1;
    z=L\(P*v); x=U\z;
    err=abs(1/m1-1/m);
    if err<=tol, break; end
    k=k+1; m=m1;
end
lam=alpha+1/m;

```

例 7.3 利用反幂法程序, 求例 7.1 中的矩阵 A 最接近 101, 99, 2 和 0 的特征值和对应的特征向量.

解 注意到此处 α 的值分别取 101, 99, 2, 0. 用反幂法的 MATLAB 程序, 编写 M 文件 ex73.m, 取容许误差为 $\varepsilon = 10^{-6}$, 在命令窗口运行该文件, 得到 4 个近似特征值: 100.7462; 99.2107; 1.7893; 0.2538, 这是矩阵 A 的两个模最大和两个模最小的特征值.

7.3 Jacobi 方法

Jacobi 方法用于求解实对称矩阵的全部特征值和对应的特征向量. 其数学原理如下:

- (1) n 阶实对称矩阵的特征值全为实数, 其对应的特征向量线性无关且两两正交.
- (2) 相似矩阵具有相同的特征值.
- (3) 若 n 阶实矩阵 A 是对称的, 则存在正交矩阵 Q , 使得 $Q^T A Q = D$, 其中 D 是一个对角矩阵, 它的对角元素 $\lambda_1, \lambda_2, \dots, \lambda_n$ 就是 A 的特征值, Q 的第 i 列向量就是 λ_i 对应的特征向量.

Jacobi 方法就是基于上述原理, 用一系列正交变换对角化 A , 即逐步消去 A 的非对角元, 从而得到 A 的全部特征值.

7.3.1 实对称矩阵的旋转正交相似变换

首先回顾一下第 2 章介绍过的一种正交变换—Givens 变换, 它是 Jacobi 方法的基本工具.

定义 7.3 设 $1 \leq i < j \leq n$, 则称矩阵

$$G_{ij} = \begin{bmatrix} 1 & & & & & \\ & \ddots & & & & \\ & & \cos \varphi & & \sin \varphi & \\ & & & 1 & & \\ & & -\sin \varphi & & \cos \varphi & \\ & & & & & \ddots \\ & & & & & & 1 \end{bmatrix} \quad (7.17)$$

为 (i, j) 平面的旋转矩阵, 或 Givens 变换矩阵.

显然, $G = G_{ij}$ 为正交矩阵, 即 $G^T G = I$. 对于向量 $x \in \mathbb{R}^n$, 由线性变换 $y = Gx$ 得到的 y 的分量为

$$\begin{cases} y_i = x_i \cos \varphi + x_j \sin \varphi, \\ y_j = -x_i \sin \varphi + x_j \cos \varphi, \\ y_k = x_k, \quad k \neq i, j, \end{cases} \quad (7.18)$$

即用 G_{ij} 对向量 x 作用, 只改变其第 i, j 两个分量.

由矩阵 $G = G_{ij}$ 确定的正交变换 $y = Gx$ 称为平面旋转变换, 或 Givens 变换. 根据式 (7.18) 容易验证, 矩阵 G_{ij} 具有下列基本性质.

定理 7.9 设 $x \in \mathbb{R}^n$ 的第 j 个分量 $x_j \neq 0$, $1 \leq i < j \leq n$. 若令

$$c = \cos \varphi = \frac{x_i}{\sqrt{x_i^2 + x_j^2}}, \quad s = \sin \varphi = \frac{x_j}{\sqrt{x_i^2 + x_j^2}}, \quad (7.19)$$

则 $y = G_{ij}x$ 的分量为

$$\begin{cases} y_i = \sqrt{x_i^2 + x_j^2}, \quad y_j = 0, \\ y_k = x_k, \quad k \neq i, j. \end{cases} \quad (7.20)$$

定理 7.9 表明, 可以用 Givens 变换将向量的某个分量变为零元素.

下面讨论 Givens 变换对实对称矩阵的作用. 用旋转矩阵 G_{ij} 对实对称矩阵 $A = (a_{ij})_{n \times n}$ 作正交相似变换, 所得矩阵记为 A_1 , 即

$$A_1 = G_{ij} A G_{ij}^T = (a_{ij}^{(1)}).$$

显然

$$A_1^T = (G_{ij} A G_{ij}^T)^T = G_{ij} A G_{ij}^T = A_1,$$

即 \mathbf{A}_1 仍为实对称矩阵. 直接计算, 得

$$\begin{cases} a_{ii}^{(1)} = a_{ii} \cos^2 \varphi + a_{jj} \sin^2 \varphi + 2a_{ij} \cos \varphi \sin \varphi, \\ a_{jj}^{(1)} = a_{ii} \sin^2 \varphi + a_{jj} \cos^2 \varphi - 2a_{ij} \cos \varphi \sin \varphi, \\ a_{ij}^{(1)} = a_{ji}^{(1)} = a_{ij}(\cos^2 \varphi - \sin^2 \varphi) - (a_{ii} - a_{jj}) \cos \varphi \sin \varphi, \\ a_{il}^{(1)} = a_{li}^{(1)} = a_{il} \cos \varphi + a_{jl} \sin \varphi, \quad l \neq i, j, \\ a_{jl}^{(1)} = a_{lj}^{(1)} = -a_{il} \sin \varphi + a_{jl} \cos \varphi, \quad l \neq i, j, \\ a_{lm}^{(1)} = a_{ml}^{(1)} = a_{ml}, \quad m, l \neq i, j. \end{cases} \quad (7.21)$$

不难看出, \mathbf{A} 经过 \mathbf{G}_{ij} 的正交相似变换后, \mathbf{A}_1 的元素和 \mathbf{A} 的元素相比, 只有第 i 行, 第 j 行和第 i 列, 第 j 列元素发生了变化, 而其他元素和 \mathbf{A} 是相同的.

由式 (7.21) 的最后一个等式可知, 若 $a_{ij} \neq 0$, 则可适当选取 φ 的值, 使得 $a_{ij}^{(1)} = a_{ji}^{(1)} = 0$. 事实上, 令

$$a_{ij}(\cos^2 \varphi - \sin^2 \varphi) - (a_{ii} - a_{jj}) \cos \varphi \sin \varphi = 0,$$

解得

$$\cot 2\varphi = \frac{a_{ii} - a_{jj}}{2a_{ij}} = \frac{1 - \tan^2 \varphi}{2 \tan \varphi}, \quad -\frac{\pi}{4} < \varphi \leq \frac{\pi}{4}. \quad (7.22)$$

在 Jacobi 方法中, 总是按上式选取 φ . 在实际计算时, 为避免使用三角函数, 可令

$$t = \tan \varphi, \quad c = \cos \varphi, \quad s = \sin \varphi, \quad d = \frac{a_{ii} - a_{jj}}{2a_{ij}}. \quad (7.23)$$

由式 (7.22), 得

$$t^2 + 2dt - 1 = 0. \quad (7.24)$$

式 (7.24) 有两个根, 取其绝对值最小者为 t , 即

$$t = \begin{cases} -d + \sqrt{d^2 + 1}, & d \geq 0, \\ -d - \sqrt{d^2 + 1}, & d < 0. \end{cases} \quad (7.25)$$

若记

$$c = \cos \varphi = \frac{1}{\sqrt{1+t^2}}, \quad s = \sin \varphi = \frac{t}{\sqrt{1+t^2}} = ct, \quad (7.26)$$

这时, 式 (7.21) 可写为

$$\begin{cases} a_{ii}^{(1)} = a_{ii}c^2 + a_{jj}s^2 + 2csa_{ij}, \\ a_{jj}^{(1)} = a_{ii}s^2 + a_{jj}c^2 - 2csa_{ij}, \quad a_{ij}^{(1)} = a_{ji}^{(1)} = 0, \\ a_{il}^{(1)} = a_{li}^{(1)} = ca_{il} + sa_{jl}, \quad a_{jl}^{(1)} = a_{lj}^{(1)} = -sa_{il} + ca_{jl}, \quad l \neq i, j; \\ a_{lm}^{(1)} = a_{ml}^{(1)} = a_{ml}, \quad m, l \neq i, j. \end{cases} \quad (7.27)$$

利用等式 $a_{ij}(c^2 - s^2) - (a_{ii} - a_{jj})cs = 0$, 可以验证

$$\left(a_{ii}^{(1)}\right)^2 + \left(a_{jj}^{(1)}\right)^2 = a_{ii}^2 + a_{jj}^2 + 2a_{ij}^2, \quad (7.28)$$

即 A 经过一次这种正交相似变换后, 所得到的矩阵 A_1 的对角元素平方和增加了 $2a_{ij}^2$.

7.3.2 Jacobi 方法及其收敛性

选择 $A_0 = A$ 中一对非零的非对角元素 a_{ij}, a_{ji} , 使用平面旋转矩阵 G_{ij} 作正交相似变换得 A_1 , 可使 A_1 的这对非对角元素 $a_{ij}^{(1)} = a_{ji}^{(1)} = 0$; 再选择 A_1 中一对非零的非对角元素作上述旋转正交相似变换得 A_2 , 可使 A_2 的这对非对角元素为 0. 如此不断地作旋转正交相似变换, 可产生一个矩阵序列 $A = A_0, A_1, \dots, A_k, \dots$. 虽然 A 至多只有 $n(n-1)/2$ 对非零非对角元素, 但不能期望通过 $n(n-1)/2$ 次旋转正交相似变换使其对角化. 因为每次旋转变换虽然能使一对特定的非对角元素化为 0, 但这次变换可能将前面已经化为 0 了的一对非对角元素变成非 0.

但是, 在 Jacobi 方法中的每一步, 如由 A_{k-1} 变成 A_k , 取其绝对值最大的一对非零非对角元素, 即取

$$|a_{ik}^{(k-1)}| = \max_{\substack{1 \leq i, j \leq n \\ i \neq j}} |a_{ij}^{(k-1)}| \quad (7.29)$$

作旋转相似变换, 这时记旋转矩阵 $G_{ij} = G_{ikjk}$. 后面将证明, 这样产生的矩阵序列 $A_0, A_1, \dots, A_k, \dots$ 趋向于对角矩阵, 即 Jacobi 方法是收敛的.

在实际计算中, 可预先取一个小的控制量 $\varepsilon > 0$, 若成立

$$|a_{ij}^{(k)}| < \varepsilon, \quad i, j = 1, 2, \dots, n, \quad i \neq j, \quad (7.30)$$

则可视 A_k 为对角矩阵, 从而结束计算. A_k 的对角元素可视作 A 的特征值.

Jacobi 方法也可以求 A 的所有特征向量. 事实上, 由

$$\begin{aligned} A_k &= G_k A_{k-1} G_k^T = G_k G_{k-1} A_{k-2} G_{k-1}^T G_k^T = \dots \\ &= G_k G_{k-1} \dots G_1 A G_1^T \dots G_{k-1}^T G_k^T, \end{aligned}$$

若记

$$Q_k = G_1^T \dots G_{k-1}^T G_k^T, \quad (7.31)$$

则

$$A_k = Q_k^T A Q_k. \quad (7.32)$$

式中: Q_k 为正交矩阵.

若 A_k 可视作对角矩阵, 其对角元即为 A 的特征值, 其第 i 个对角元 $a_{ii}^{(k)}$ 对应的特征向量就是 Q_k 的第 i 列元素构成的向量. Q_k 的计算可与 A 的旋转相似变换同步进行. 若令 $Q_0 = I$, 则

$$Q_k = Q_{k-1} G_k^T. \quad (7.33)$$

若 $G_k = G_{ij}$, 得 Q_k 的计算公式如下:

$$\begin{cases} q_{li}^{(k)} = q_{li}^{(k-1)}c + q_{lj}^{(k-1)}s, & l = 1, 2, \dots, n, \\ q_{lj}^{(k)} = -q_{li}^{(k-1)}s + q_{lj}^{(k-1)}c, & l = 1, 2, \dots, n, \\ q_{km}^{(k)} = q_{km}^{(k-1)}, & k, m \neq i, j. \end{cases} \quad (7.34)$$

也就是说, 除了第 i, j 列元素发生变化外, 其他元素不变. 若不需要计算特征向量, 则可省略此步.

根据上述讨论, 可得 Jacobi 方法的计算步骤如下.

算法 7.3 (Jacobi 方法)

步 1, 输入矩阵 A , $Q = I$, 初始向量 x , 误差限 ε , 最大迭代次数 N . 置 $k := 1$.

步 2, 在矩阵中找绝对值最大的非对角元

$$\mu = |a_{i_r, j_r}| = \max_{\substack{1 \leq i, j \leq n \\ i \neq j}} |a_{ij}|,$$

置 $i := i_r, j := j_r$.

步 3, 按式 (7.23) ~ 式 (7.27) 计算 d, t, c, s 的值和矩阵 A_1 的元素 $a_{lm}^{(1)}$, $l, m = 1, 2, \dots, n$.

步 4, 更新 Q 的元素:

$$\begin{cases} q_{li} := q_{li}c + q_{lj}s, \\ q_{lj} := -q_{li}s + q_{lj}c, \end{cases} \quad l = 1, 2, \dots, n.$$

步 5, 若 $\mu < \varepsilon$, 输出 A_1 的对角元和 Q 的列向量, 停算; 否则, 转步 6.

步 6, 若 $k < N$, 置 $k := k + 1$, 转步 2; 否则输出计算失败信息, 停算.

根据算法 7.3, 编制 Jacobi 方法的 MATLAB 程序如下:

```
%Jacobi方法程序-Jacobi_eig.m
function [lambda,Q]=Jacobi_eig(A,tol)
%用Jacobi方法求实对称矩阵A的全部特征值和特征向量
%输入:A为n阶对称方阵,tol为容许误差
%输出:lambda为向量,其分量为A的特征值,
%      Q为矩阵,其元素为矩阵A的n个特征向量
if nargin<2, tol=1e-6; end
[n]=size(A,1); Q=eye(n);
%计算A的非对角元绝对值最大元素所在的行p和列q
[w1,p]=max(abs(A-diag(diag(A)))));
[w2,q]=max(w1); p=p(q);
while(1)
    d=(A(p,p)-A(q,q))/(2*A(p,q));
```

```

if(d>=0)
    t=-d+sqrt(d^2+1);
else
    t=-d-sqrt(d^2+1);
end
c=1/sqrt(t^2+1); s=c*t; G=[c s; -s c];
A([p q],:)=G*A([p q],:);
A(:,[p q])=A(:,[p q])*G';
Q(:,[p q])=Q(:,[p q])*G';
[w1,p]=max(abs(A-diag(diag(A))));
[w2,q]=max(w1); p=p(q);
if (abs(A(p,q))<tol*sqrt(sum(diag(A).^2)/n))
    break;
end
end
lambda=sort(diag(A));

```

例 7.4 利用 Jacobi 方法程序, 求例 7.1 中的矩阵 A 的全部特征值和对应的特征向量.

解 在 MATLAB 命令窗口执行程序 ex74.m 得到矩阵 A 的全部特征值 $\hat{\lambda}$. 此外, 利用 MATLAB 自带的函数 eig 求得其特征值为 λ , 计算其误差 $\|\hat{\lambda} - \lambda\|_2 = 3.1756 \times 10^{-9}$. 特征值的分布如图 7.1 所示.

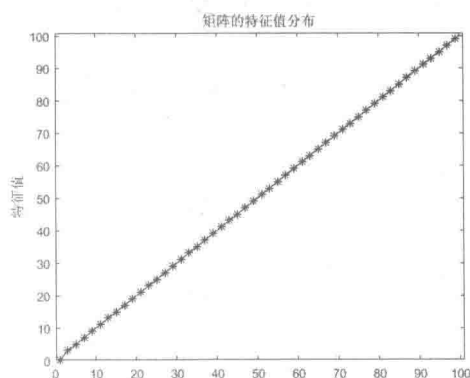


图 7.1 矩阵 A 的特征值分布

下面考虑算法 7.3 (Jacobi 方法) 的收敛性. 记实对称矩阵 A 的非对角元素的平方和为

$$S(A) = \sum_{\substack{l,m=1 \\ l \neq m}}^n a_{lm}^2. \quad (7.35)$$

设 $A_{k+1} = G_{ij} A_k G_{ij}^T$, 则由式 (7.27) 不难验证

$$S(A_{k+1}) = S(A_k) - 2[a_{ij}^{(k)}]^2. \quad (7.36)$$

即经过这种正交相似变换后, A_{k+1} 的非对角元素平方和减少了 $2[a_{ij}^{(k)}]^2$. 同时, 由式 (7.28), 有

$$[a_{ii}^{(k+1)}]^2 + [a_{jj}^{(k+1)}]^2 = [a_{ii}^{(k)}]^2 + [a_{jj}^{(k)}]^2 + 2[a_{ij}^{(k)}]^2, \quad (7.37)$$

即对角元素的平方和增加了 $2[a_{ij}^{(k)}]^2$. 若在 Jacobi 方法中, 每次旋转正交相似变换使 A_k 的绝对值最大的非对角元素化为 0, 则成立以下定理.

定理 7.10 记实对称矩阵 $A = A_0$, 若在 Jacobi 方法中, 每次旋转正交相似变换使 A_k 的绝对值最大的非对角元素化为 0, 则得到的矩阵序列 $\{A_k\}$ 趋向于对角矩阵.

证明 设 A_k 的绝对值最大的非对角元素为 $a_{ij}^{(k)}$, 故有

$$[a_{ij}^{(k)}]^2 \geq \frac{1}{n(n-1)} S(A_k).$$

用旋转正交相似变换将其化为 0, 得 A_{k+1} , 此时

$$\begin{aligned} S(A_{k+1}) &= S(A_k) - 2[a_{ij}^{(k)}]^2 \leq S(A_k) - \frac{2}{n(n-1)} S(A_k) \\ &= \left[1 - \frac{2}{n(n-1)}\right] S(A_k) \leq \left[1 - \frac{2}{n(n-1)}\right]^{k+1} S(A). \end{aligned}$$

由于

$$1 - \frac{2}{n(n-1)} < 1,$$

所以

$$\lim_{k \rightarrow \infty} S(A_{k+1}) = 0,$$

即 A_{k+1} 趋向于对角矩阵, 故 Jacobi 方法是收敛的. 证毕. \square

7.4 QR 方法

QR 方法用于求一般矩阵的全部特征值, 是目前最有效的方法之一. 本节就实矩阵的情形进行介绍.

众所周知, 对于任何实对称矩阵 $A \in \mathbb{R}^{n \times n}$, 存在正交矩阵 Q 使得

$$Q^T A Q = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n),$$

式中: $\lambda_1, \lambda_2, \dots, \lambda_n$ 为 A 的全部特征值; Q 的列向量为对应的特征向量.

而对于一般的 $A \in \mathbb{R}^{n \times n}$, 有下面的实 Schur 分解定理.

定理 7.11 对于任何实矩阵 $A \in \mathbb{R}^{n \times n}$, 存在正交矩阵 Q 使得

$$Q^T A Q = \begin{bmatrix} R_{11} & R_{12} & \cdots & R_{1m} \\ & R_{22} & \cdots & R_{2m} \\ & & \ddots & \vdots \\ & & & R_{mm} \end{bmatrix}, \quad (7.38)$$

式中: 对角块 $R_{ii} (i = 1, 2, \dots, m)$ 为 1×1 或 2×2 的子矩阵, 1 阶子矩阵的元素是 A 的实特征值, 2 阶子矩阵的两个特征值是 A 的一对复共轭特征值.

式 (7.38) 通常称为矩阵 A 的实 Schur 分解, 而右边的分块上三角矩阵称为 A 的实 Schur 标准形. 显然, 只要求得一个实矩阵的实 Schur 标准形, 就很容易求得它的全部特征值.

因此, 通常希望构造一种迭代 (相似变换), 希望它能逼近矩阵 A 的实 Schur 标准形. 例如, 对于给定的矩阵 $A \in \mathbb{R}^{n \times n}$, 令 $A_1 := A$, 构造迭代:

$$\begin{cases} A_k = Q_k R_k, \\ A_{k+1} = R_k Q_k, \end{cases} \quad k = 1, 2, \dots, \quad (7.39)$$

式中: Q_k 为正交矩阵; R_k 为上三角矩阵.

可以证明, 在一定条件下, 由式 (7.39) 产生的矩阵序列 $\{A_k\}$ 将“逼近”于 A 的实 Schur 标准形.

然而, 式 (7.39) 作为一种实用的迭代法是没有竞争力的. 一是每步迭代的运算量太大 (大约是 $O(n^3)$); 二是收敛速度太缓慢 (依赖于特征值的分离程度). 因此, 要想其成为一种高效的迭代方法, 必须设法尽可能地减少每步迭代的运算量, 提高其收敛速度. 一个可行的办法是, 首先把矩阵 A 正交相似变换为上 Hessenberg 形, 然后再用正交相似变换对它进行迭代. 下面就循着这个思路介绍实用的 QR 方法来求实矩阵 A 的全部特征值.

7.4.1 化一般矩阵为上 Hessenberg 矩阵

在用 QR 方法求矩阵特征值时, Householder 矩阵有两个作用: 一是对 A 作正交相似变换, 把 A 化为上 Hessenberg 矩阵; 二是对矩阵作正交三角分解.

首先讨论把 A 化为上 Hessenberg 矩阵. 设 $A_1 = A = (a_{ij}^{(1)})$ 是 n 阶实方阵, 取 $x = (0, a_{21}^{(1)}, \dots, a_{n1}^{(1)})^T$, 记 $a_1 = \text{sgn}(x_2) \|x\|_2$, 则由定理 2.2 和定理 2.3 构造 Householder 矩阵

$$H_1 = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & * & \cdots & * \\ \vdots & \vdots & \ddots & \vdots \\ 0 & * & \cdots & * \end{bmatrix},$$

使得

$$H_1 x = a_1 e_2.$$

所以 $H_1 A_1$ 的第 1 列为

$$H_1 \begin{bmatrix} a_{11}^{(1)} \\ a_{21}^{(1)} \\ \vdots \\ a_{n1}^{(1)} \end{bmatrix} = H_1 x + H_1 \begin{bmatrix} a_{11}^{(1)} \\ 0 \\ \vdots \\ 0 \end{bmatrix} = a_1 e_2 + \begin{bmatrix} a_{11}^{(1)} \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \begin{bmatrix} a_{11}^{(1)} \\ a_1 \\ \vdots \\ 0 \end{bmatrix}.$$

因为用 H_1 右乘一个矩阵不改变该矩阵的第 1 列, 于是

$$A_2 = H_1 A_1 H_1 = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(2)} & \cdots & a_{1n}^{(2)} \\ a_1 & a_{22}^{(2)} & \cdots & a_{2n}^{(2)} \\ 0 & a_{32}^{(2)} & \cdots & a_{3n}^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & a_{n2}^{(2)} & \cdots & a_{nn}^{(2)} \end{bmatrix}.$$

再取 $x = (0, 0, a_{32}^{(2)}, \cdots, a_{n2}^{(2)})^T$, 记 $a_2 = \operatorname{sgn}(x_3) \|x\|_2$, 构造 H_2 为

$$H_2 = \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & * & \cdots & * \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & * & \cdots & * \end{bmatrix},$$

使得

$$H_2 x = a_2 e_3.$$

所以 $H_2 A_2$ 的第 1 列与 A_2 的第 1 列相同, 而 $H_2 A_2$ 的第 2 列变为

$$H_2 x + H_2 \begin{bmatrix} a_{12}^{(2)} \\ a_{22}^{(2)} \\ 0 \\ \vdots \\ 0 \end{bmatrix} = a_2 e_3 + \begin{bmatrix} a_{12}^{(2)} \\ a_{22}^{(2)} \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \begin{bmatrix} a_{12}^{(2)} \\ a_{22}^{(2)} \\ a_2 \\ \vdots \\ 0 \end{bmatrix}.$$

而用 H_2 右乘一个矩阵不改变该矩阵的第 1 列和第 2 列, 于是

$$A_3 = H_2 A_2 H_2 = \begin{bmatrix} * & * & * & \cdots & * \\ a_1 & * & * & \cdots & * \\ 0 & a_2 & * & \cdots & * \\ 0 & 0 & * & \cdots & * \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & * & \cdots & * \end{bmatrix}.$$

这样下去, 经过 $n-2$ 次变换后, A_1 就化为上 Hessenberg 矩阵 A_{n-1} , 即

$$A_{n-1} = H_{n-2} \cdots H_2 H_1 A_1 H_1 H_2 \cdots H_{n-2}$$

$$= \begin{bmatrix} * & * & * & * & \cdots & * \\ a_1 & * & * & * & \cdots & * \\ & a_2 & * & * & \cdots & * \\ & & a_3 & * & \cdots & * \\ & & & \ddots & \ddots & \vdots \\ & & & & a_{n-1} & * \end{bmatrix},$$

如果 A_1 是对称矩阵, 则 A_{n-1} 仍是对称矩阵, 此时 A_{n-1} 将是对称三对角矩阵:

$$A_{n-1} = \begin{bmatrix} * & a_1 & & & \\ a_1 & * & a_2 & & \\ & a_2 & * & a_3 & \\ & & \ddots & \ddots & a_{n-1} \\ & & & a_{n-1} & * \end{bmatrix}.$$

以上利用 Householder 变换约化 A 为上 Hessenberg 形的方法, 可总结为如下实用算法.

- 算法 7.4 (上 Hessenberg 化)**
- 步 1, 输入 $A = (a_{ij})$, $k := 1$.
 - 步 2, 计算 $n-k$ 阶 Householder 矩阵 \widetilde{H}_k , 使

$$\widetilde{H}_k \begin{bmatrix} a_{k+1,k} \\ a_{k+2,k} \\ \vdots \\ a_{n,k} \end{bmatrix} = \begin{bmatrix} * \\ 0 \\ \vdots \\ 0 \end{bmatrix},$$

- 置 $A := H_k A H_k$, 其中 $H_k = \text{diag}(I_k, \widetilde{H}_k)$.
- 步 3, 若 $k < n-2$, 则 $k := k+1$, 转步 2; 否则输出有关信息, 停算.

算法 7.4 计算出 A 的上 Hessenberg 形就存在 A 所对应的存储单元内, 所需运算量是 $5n^3/3$. 如果需要累计 $Q = H_1 H_2 \cdots H_{n-2}$, 则还需增加 $2n^3/3$ 的运算量.

将矩阵 A 化为上 Hessenberg 矩阵 (算法 7.4) 的 MATLAB 程序如下:

```
function [A,Q]=mhessen(A)
%用Householder变换化n阶矩阵A为上Hessenberg矩阵
%调用函数:mhouse.m
```

```

n=size(A,1); Q=eye(n);
for k=1:(n-2)
    x=A(k+1:n,k); [v,beta]=mhouse(x);
    H=(eye(length(v))-beta*v*v');
    A(k+1:n,1:n)=H*A(k+1:n,1:n);
    A(1:n,k+1:n)=A(1:n,k+1:n)*H;
    Q=Q*blkdiag(eye(k),H);
end

```

例 7.5 利用 MATLAB 程序 mhessen.m, 将下列矩阵化为上 Hessenberg 矩阵:

$$A = \begin{bmatrix} -1 & 2 & 3 & 5 \\ 2 & -3 & 8 & 1 \\ 3 & 8 & -2 & 7 \\ 5 & 1 & 7 & 6 \end{bmatrix}.$$

解 在 MATLAB 命令窗口输入:

```
>> A=[-1 2 3 5; 2 -3 8 1; 3 8 -2 7; 5 1 7 6];
```

```
>> [A,Q]=mhessen(A)
```

A =

```

-1.0000    -6.1644     0.0000     0.0000
-6.1644    11.7368     1.8380     0.0000
 0.0000     1.8380    -6.5929     5.9938
 0.0000     0.0000     5.9938    -4.1439

```

Q =

```

1.0000         0         0         0
         0     0.3244    -0.0418    -0.9450
         0     0.4867     0.8640     0.1289
         0     0.8111    -0.5017     0.3007

```

一般来说, 上 Hessenberg 分解是不唯一的, 但可以证明下面的结果.

定理 7.12 设 $A \in \mathbb{R}^{n \times n}$ 有如下两个上 Hessenberg 分解:

$$U^T A U = H, \quad V^T A V = G, \quad (7.40)$$

式中: $U = [u_1, u_2, \dots, u_n]$ 和 $V = [v_1, v_2, \dots, v_n]$ 为 n 阶正交矩阵; $H = (h_{ij})$ 和 $G = (g_{ij})$ 为上 Hessenberg 矩阵. 若 $u_1 = v_1$, 且 H 的次对角元 $h_{i+1,i} \neq 0$ ($i = 1, 2, \dots, n$), 则存在对角元均为 1 或 -1 的对角矩阵 D , 使得

$$U = V D, \quad H = D G D, \quad (7.41)$$

即 $u_i = \pm v_i$, $|h_{ij}| = |g_{ij}|$, $i, j = 1, 2, \dots, n$.

证明 对矩阵的阶数 n 用归纳法. 对于 $n = 1$ 时结论显然成立. 假设 $n = m$ 时结论成立, 即

$$\mathbf{u}_i = \varepsilon_i \mathbf{v}_i, \quad i = 1, 2, \dots, m, \quad (7.42)$$

式中: $\varepsilon_1 = 1$, $\varepsilon_i = 1$ 或 -1 , $i = 2, \dots, m$. 下面证明存在 ε_{m+1} 为 1 或 -1 使得

$$\mathbf{u}_{m+1} = \varepsilon_{m+1} \mathbf{v}_{m+1}.$$

由式 (7.40), 得

$$\mathbf{A}\mathbf{U} = \mathbf{U}\mathbf{H}, \quad \mathbf{A}\mathbf{V} = \mathbf{V}\mathbf{G}.$$

分别比较上面两个矩阵等式两边的第 m 列, 可得

$$\mathbf{A}\mathbf{u}_m = h_{1m}\mathbf{u}_1 + \dots + h_{mm}\mathbf{u}_m + h_{m+1,m}\mathbf{u}_{m+1}, \quad (7.43)$$

$$\mathbf{A}\mathbf{v}_m = g_{1m}\mathbf{v}_1 + \dots + g_{mm}\mathbf{v}_m + g_{m+1,m}\mathbf{v}_{m+1}. \quad (7.44)$$

分别在式 (7.43) 和式 (7.44) 两边左乘 \mathbf{u}_i^T 和 \mathbf{v}_i^T ($i = 1, 2, \dots, m$), 得

$$h_{im} = \mathbf{u}_i^T \mathbf{A}\mathbf{u}_m, \quad i = 1, 2, \dots, m, \quad (7.45)$$

$$g_{im} = \mathbf{v}_i^T \mathbf{A}\mathbf{v}_m, \quad i = 1, 2, \dots, m. \quad (7.46)$$

由式 (7.42)、式 (7.45) 和式 (7.46), 得

$$h_{im} = \varepsilon_i \varepsilon_m g_{im}, \quad i = 1, 2, \dots, m. \quad (7.47)$$

将式 (7.47) 代入式 (7.43), 并利用式 (7.42) 和式 (7.44), 得

$$\begin{aligned} h_{m+1,m}\mathbf{u}_{m+1} &= \mathbf{A}\mathbf{u}_m - \varepsilon_1 \varepsilon_m g_{1m}\mathbf{u}_1 - \dots - \varepsilon_m \varepsilon_m g_{mm}\mathbf{u}_m \\ &= \varepsilon_m (\mathbf{A}\mathbf{v}_m - \varepsilon_1^2 g_{1m}\mathbf{v}_1 - \dots - \varepsilon_m^2 g_{mm}\mathbf{v}_m) \\ &= \varepsilon_m (\mathbf{A}\mathbf{v}_m - g_{1m}\mathbf{v}_1 - \dots - g_{mm}\mathbf{v}_m) \\ &= \varepsilon_m g_{m+1,m}\mathbf{v}_{m+1}. \end{aligned} \quad (7.48)$$

上式两边取范数, 得

$$|h_{m+1,m}| = |g_{m+1,m}|.$$

由于 $h_{m+1,m} \neq 0$, 故式 (7.48) 蕴含着

$$\mathbf{u}_{m+1} = \varepsilon_{m+1} \mathbf{v}_{m+1},$$

式中: $\varepsilon_{m+1} = 1$ 或 -1 . 证毕. □

注 7.3 一个上 Hessenberg 矩阵 $\mathbf{H} = (h_{ij})$, 如果其次对角元均不为零, 即 $h_{i+1,i} \neq 0$, $i = 1, 2, \dots, n-1$, 则它是不可约的. 定理 7.12 表明, 如果 $\mathbf{Q}^T \mathbf{A} \mathbf{Q} = \mathbf{H}$ 为不可约的上 Hessenberg 矩阵, 则 \mathbf{Q} 和 \mathbf{H} 完全由 \mathbf{Q} 的第 1 列确定 (这里是在相差一个正负号意义下的唯一).

7.4.2 上 Hessenberg 矩阵的 QR 分解

对于上 Hessenberg 矩阵

$$H = \begin{bmatrix} h_{11}^{(1)} & h_{12}^{(1)} & h_{13}^{(1)} & \cdots & h_{1n}^{(1)} \\ h_{21}^{(1)} & h_{22}^{(1)} & h_{23}^{(1)} & \cdots & h_{2n}^{(1)} \\ & h_{32}^{(1)} & h_{33}^{(1)} & \cdots & h_{3n}^{(1)} \\ & & \ddots & \ddots & \vdots \\ & & & h_{n,n-1}^{(1)} & h_{nn}^{(1)} \end{bmatrix},$$

通常可以通过 $n-1$ 次 Givens 变换将它化成上三角矩阵, 从而得到 H 的 QR 分解式. 具体步骤是:

(1) 记 $H_1 = H$. 设 $h_{21}^{(1)} \neq 0$ (否则可进行下一步), 取 Givens 矩阵

$$G_{21} = \begin{bmatrix} c_1 & s_1 & & & \\ -s_1 & c_1 & & & \\ & & 1 & & \\ & & & \ddots & \\ & & & & 1 \end{bmatrix},$$

式中:

$$c_1 = \frac{h_{11}^{(1)}}{r_1}, \quad s_1 = \frac{h_{21}^{(1)}}{r_1}, \quad r_1 = \sqrt{(h_{11}^{(1)})^2 + (h_{21}^{(1)})^2}.$$

则

$$G_{21}H_1 = \begin{bmatrix} r_1 & h_{12}^{(2)} & h_{13}^{(2)} & \cdots & h_{1n}^{(2)} \\ 0 & h_{22}^{(2)} & h_{23}^{(2)} & \cdots & h_{2n}^{(2)} \\ & h_{32}^{(2)} & h_{33}^{(2)} & \cdots & h_{3n}^{(2)} \\ & & \ddots & \ddots & \vdots \\ & & & h_{n,n-1}^{(2)} & h_{nn}^{(2)} \end{bmatrix} := H_2.$$

(2) 设 $h_{32}^{(1)} \neq 0$ (否则可进行下一步), 再取 Givens 矩阵

$$G_{32} = \begin{bmatrix} 1 & & & & \\ & c_2 & s_2 & & \\ & -s_2 & c_2 & & \\ & & & 1 & \\ & & & & \ddots \\ & & & & & 1 \end{bmatrix},$$

式中:

$$c_2 = \frac{h_{22}^{(2)}}{r_2}, \quad s_2 = \frac{h_{32}^{(2)}}{r_2}, \quad r_2 = \sqrt{(h_{22}^{(2)})^2 + (h_{32}^{(2)})^2}.$$

则

$$G_{32}H_2 = \begin{bmatrix} r_1 & h_{12}^{(3)} & h_{13}^{(3)} & \cdots & h_{1,n-1}^{(3)} & h_{1n}^{(3)} \\ 0 & r_2 & h_{23}^{(3)} & \cdots & h_{2,n-1}^{(3)} & h_{2n}^{(3)} \\ & 0 & h_{33}^{(3)} & \cdots & h_{3,n-1}^{(3)} & h_{3n}^{(3)} \\ & & h_{43}^{(3)} & \cdots & h_{4,n-1}^{(3)} & h_{4n}^{(3)} \\ & & & \ddots & \vdots & \vdots \\ & & & & h_{n,n-1}^{(3)} & h_{nn}^{(3)} \end{bmatrix} := H_3.$$

(3) 假设上述过程已经进行了 $k-1$ 步, 有

$$H_k = G_{k,k-1}H_{k-1} = \begin{bmatrix} r_1 & \cdots & h_{1,k-1}^{(k)} & h_{1k}^{(k)} & \cdots & h_{1,n-1}^{(k)} & h_{1n}^{(k)} \\ & \ddots & & & & & \\ & & r_{k-1} & h_{k-1,k}^{(k)} & \cdots & h_{k-1,n-1}^{(k)} & h_{k-1,n}^{(k)} \\ & & & h_{kk}^{(k)} & \cdots & h_{k,n-1}^{(k)} & h_{kn}^{(k)} \\ & & & h_{k+1,k}^{(k)} & \cdots & h_{k+1,n-1}^{(k)} & h_{k+1,n}^{(k)} \\ & & & & \ddots & \vdots & \vdots \\ & & & & & h_{n,n-1}^{(k)} & h_{nn}^{(k)} \end{bmatrix}.$$

设 $h_{k+1,k}^{(k)} \neq 0$, 取 Givens 矩阵

$$G_{k+1,k} = \begin{bmatrix} 1 & & & & & & \\ & \ddots & & & & & \\ & & 1 & & & & \\ & & & c_k & s_k & & \\ & & & -s_k & c_k & & \\ & & & & & 1 & \\ & & & & & & \ddots \\ & & & & & & & 1 \end{bmatrix},$$

式中:

$$c_k = \frac{h_{kk}^{(k)}}{r_k}, \quad s_k = \frac{h_{k+1,k}^{(k)}}{r_k}, \quad r_k = \sqrt{(h_{kk}^{(k)})^2 + (h_{k+1,k}^{(k)})^2}.$$

于是

$$G_{k+1,k}H_k = \begin{bmatrix} r_1 & \cdots & h_{1,k}^{(k+1)} & h_{1,k+1}^{(k+1)} & \cdots & h_{1,n-1}^{(k+1)} & h_{1n}^{(k+1)} \\ & & & & & & \\ & & & & & & \\ & & & r_k & h_{k,k+1}^{(k+1)} & \cdots & h_{k,n-1}^{(k+1)} & h_{kn}^{(k+1)} \\ & & & h_{k+1,k+1}^{(k+1)} & \cdots & h_{k+1,n-1}^{(k+1)} & h_{k+1,n}^{(k+1)} \\ & & & h_{k+2,k+1}^{(k+1)} & \cdots & h_{k+2,n-1}^{(k+1)} & h_{k+2,n}^{(k+1)} \\ & & & & \ddots & \vdots & \vdots \\ & & & & & h_{n,n-1}^{(k+1)} & h_{nn}^{(k+1)} \end{bmatrix} := H_{k+1}.$$

因此,最多作 $n-1$ 次 Givens 变换,即得

$$G_{n,n-1} \cdots G_{32} G_{21} H = \begin{bmatrix} r_1 & h_{12}^{(n)} & h_{13}^{(n)} & \cdots & h_{1n}^{(n)} \\ & r_2 & h_{23}^{(n)} & \cdots & h_{2n}^{(n)} \\ & & r_3 & \cdots & h_{3n}^{(n)} \\ & & & \ddots & \vdots \\ & & & & r_n \end{bmatrix} = R.$$

因为 $G_{k,k-1}$ ($k=2, \cdots, n$) 均为正交阵,故

$$H = G_{21}^T G_{32}^T \cdots G_{n,n-1}^T R = QR,$$

式中: $Q = G_{21}^T G_{32}^T \cdots G_{n,n-1}^T$ 仍为正交阵.

可算出完成这一过程的运算量约为 $4n^2$, 比一般矩阵的 QR 分解的运算量 $O(n^3)$ 少了一个数量级.

值得注意的是,可以证明 $\widetilde{H} = RQ = Q^T H Q$ 仍为上 Hessenberg 矩阵,于是可按上述步骤一直迭代下去,直到 H 正交相似于上三角矩阵或块上三角矩阵 (对角块为 1×1 或 2×2 矩阵) 为止,从而求得矩阵 H 的全部特征值和相应的特征向量.

上 Hessenberg 矩阵 QR 分解的 MATLAB 程序如下:

```
function A=hessen_qrtran(A,m)
%本程序输入n阶上Hessenberg矩阵A,用Givens变换对其左上角m阶主子块
%进行QR分解,再作相似变换,最后输出变换后的上Hessenberg矩阵A.
Q=eye(m);
for i=1:m-1
    xi=A(i,i); xk=A(i+1,i);
    if xk~=0
        d=sqrt(xi^2+xk^2);
```

```

        c=xi/d; s=xk/d;
        G=[c, s; -s, c];
        A(i:i+1,i:m)=G*A(i:i+1,i:m);
        Q(1:m,i:i+1)=Q(1:m,i:i+1)*G';
    end
end
%Q*A, %验证Q*R=A
A(1:m,1:m)=A(1:m,1:m)*Q;

```

例 7.6 利用程序将上 Hessenberg 矩阵 A 进行 QR 变换, 其中

$$A = \begin{bmatrix} 2 & 3 & 5 & 7 & 8 \\ 4 & 2 & 3 & 5 & 9 \\ 0 & 8 & 3 & 6 & 2 \\ 0 & 0 & 7 & 1 & 3 \\ 0 & 0 & 0 & 6 & 9 \end{bmatrix}.$$

解 在 MATLAB 命令窗口输入:

```

>> A=[2 3 5 7 8; 4 2 3 5 9; 0 8 3 6 2; 0 0 7 1 3; 0 0 0 6 9];
>> A=hessen_qrtran(A,5)
A =

```

```

    4.8000    5.3682    7.6664   10.5793   -4.9445
    7.3321    2.7238    4.9535    0.3480   -6.8327
         0    7.2218    1.1745    2.7581   -3.2311
         0         0    6.0298    7.7534   -4.9785
         0         0         0   -1.5283    0.5483

```

7.4.3 基本 QR 方法

本节介绍求一般方阵全部特征值的 QR 方法. 令 $A_1 = A$, 对 A_1 作 QR 分解:

$$A_1 = Q_1 R_1,$$

然后令 $A_2 = R_1 Q_1$, 再对 A_2 作 QR 分解:

$$A_2 = Q_2 R_2,$$

并令 $A_3 = R_2 Q_2$, 这样下去就得到一个矩阵序列 $\{A_k\}$, 其产生过程可概述如下:

$$\begin{cases} A_1 = A, \\ A_k = Q_k R_k, \\ A_{k+1} = R_k Q_k, \end{cases} \quad k = 1, 2, \dots. \quad (7.49)$$

容易证明, A_{k+1} 与 A_k 相似, 故 $\{A_k\}$ 有相同的特征值.

在一定条件下, $\{A_k\}$ 本质上收敛于上三角矩阵 (或分块上三角矩阵). 若它们收敛于上三角矩阵, 则该上三角矩阵的对角元就是原矩阵 A 的全部特征值; 若收敛于分块上三角矩阵, 则这些分块矩阵的特征值也就是 A 的特征值.

由于当 A 为一般的实矩阵时, $\{A_k\}$ 的收敛速度较慢, 故在 QR 方法的实际应用中, 通常先将 A 化为相似的上 Hessenberg 矩阵, 再求特征值以加快收敛速度. 它的计算过程如下.

算法 7.5 (基本 QR 方法)

步 1, 输入上 Hessenberg 矩阵 $A \in \mathbb{R}^{n \times n}$.

步 2, 记 $A_1 := A$. 对于 $k = 1, 2, \dots$, 有

$$(1) A_k = Q_k R_k \quad (\text{QR 分解}).$$

$$(2) A_{k+1} = Q_k^T A_k Q_k = R_k Q_k \quad (\text{正交相似变换}).$$

基本 QR 方法的 MATLAB 程序如下:

```
function [iter,D]=qr_eig(A,tol,N)
%用基本QR算法求n阶实方阵A的全部特征值
%输入:A为实对称矩阵,tol为控制精度,N为最大迭代次数
%输出:iter为迭代次数,D为A的全部特征值
%调用函数:mhessen.m,hessen_qrtran.m,eig-仅用于1,2矩阵
if nargin<3, N=500; end
if nargin<2, tol=1e-5; end
n=size(A,1); D=zeros(n,1);
i=n; m=n; iter=0; %初始化
A=mhessen(A); %化矩阵A为Hessenberg矩阵
while (iter<=N) %用基本QR算法进行迭代
    iter=iter+1;
    if m<=2
        la=eig(A(1:m,1:m)); D(1:m)=la';
        break;
    end
    %对上Hessenberg矩阵作QR分解并作正交相似变换
    A=hessen_qrtran(A,m);
    %下面的程序段判断是否终止
    for k=m-1:-1:1
        if abs(A(k+1,k))<tol
            if m-k<=2
                la=eig(A(k+1:m,k+1:m));
                j=i-m+k+1; D(j:i)=la';
                i=j-1; m=k; break;
            end
        end
    end
end
```



```

        end
    end
end
end

```

例 7.7 利用程序 qr_eig.m, 求下列矩阵的全部特征值:

$$A = \begin{bmatrix} 3 & 2 & 3 & 4 & 5 & 6 & 7 \\ 11 & 1 & 2 & 3 & 4 & 5 & 6 \\ 2 & 8 & 9 & 1 & 2 & 3 & 4 \\ -4 & 2 & 9 & 11 & 13 & 15 & 8 \\ -1 & -2 & -3 & -1 & -1 & -1 & -1 \\ 3 & 2 & 3 & 4 & 13 & 15 & 8 \\ -2 & -2 & -3 & -4 & -5 & -3 & -3 \end{bmatrix}.$$

解 在 MATLAB 命令窗口输入:

```

>> A=[3 2 3 4 5 6 7;11 1 2 3 4 5 6;2 8 9 1 2 3 4; ...
      -4 2 9 11 13 15 8; -1 -2 -3 -1 -1 -1 -1; ...
      3 2 3 4 13 15 8; -2 -2 -3 -4 -5 -3 -3];
>> [iter,D]=qr_eig(A)
iter =
      622
D =
 18.4123 + 0.0000i
 11.1805 + 0.0000i
  1.7099 - 4.2522i
  1.7099 + 4.2522i
  4.4983 + 0.0000i
 -2.2327 + 0.0000i
 -0.2783 + 0.0000i

```

下面分析基本 QR 方法的收敛性.

定义 7.4 若由 QR 方法产生的序列 $\{A_k\}$ 当 $k \rightarrow \infty$ 收敛于分块上三角矩阵 (对角块为一阶或二阶子块), 则称 QR 方法是收敛的. 若序列 $\{A_k\}$ 当 $k \rightarrow \infty$ 时, 其对角元均收敛且严格下三角部分元素收敛于 0, 则称 $\{A_k\}$ 基本收敛到上三角阵.

值得注意的是, 基本收敛的概念并未指出 $\{A_k\}$ 严格上三角部分元素是否收敛. 但对求矩阵 A 的特征值而言, 基本收敛足够了.

算法 7.5 具有下列性质.

性质 7.1 在算法 7.5 中, 若记

$$\tilde{Q}_k = Q_1 Q_2 \cdots Q_k, \quad \tilde{R}_k = R_k R_{k-1} \cdots R_1, \quad (7.50)$$

显然 \tilde{Q}_k 为正交矩阵, \tilde{R}_k 为上三角矩阵. 则有

(1) Q_k 和 A_{k+1} 都是上 Hessenberg 矩阵.

$$(2) A_{k+1} = \tilde{Q}_k^T A \tilde{Q}_k \sim A \quad (\text{相似性}). \quad (7.51)$$

$$(3) A^k = \tilde{Q}_k \tilde{R}_k \quad (A \text{ 的 } k \text{ 次幂 } A^k \text{ 的 QR 分解}). \quad (7.52)$$

证明 用归纳法证明性质 (3). 当 $k=1$ 时, 有

$$A = A_1 = Q_1 R_1 = \tilde{Q}_1 \tilde{R}_1.$$

设 $A^{k-1} = \tilde{Q}_{k-1} \tilde{R}_{k-1}$, 则

$$\begin{aligned} A^k &= A(\tilde{Q}_{k-1} \tilde{R}_{k-1}) = \tilde{Q}_{k-1} (\tilde{Q}_{k-1}^T A \tilde{Q}_{k-1}) \tilde{R}_{k-1} \\ &= \tilde{Q}_{k-1} A_k \tilde{R}_{k-1} = \tilde{Q}_{k-1} Q_k R_k \tilde{R}_{k-1} = \tilde{Q}_k \tilde{R}_k. \end{aligned}$$

证毕. □

性质 (1) 称为上 Hessenberg 形在 QR 变换下的不变性. 它的意义是, 算法 7.5 可以始终对上 Hessenberg 矩阵进行. 这时, QR 分解中每列的消元只要作一次 Givens 变换, 从而简化了 QR 变换的计算. 性质 (2) 表示, 由 QR 方法生成的矩阵序列 $\{A_k\}$ 保持原矩阵 A 的特征值不变. 性质 (3) 说明, QR 方法即 QR 变换过程, 实质上是对 A 的 k 次幂 A^k 进行 QR 分解的过程. 由此可见, QR 方法与幂法有内在的联系. 下面给出 QR 方法的一个最简单的收敛性定理, 它表明, 在一定条件下可以把 QR 方法看成幂法的推广. 为此, 先给出下面的引理.

引理 7.1 设 $Q = [q_1, Q_{n-1}] \in \mathbb{R}^{n \times n}$ 为正交矩阵, 其中 q_1 为 Q 的第 1 列. 对于矩阵 A , 若记

$$Q^T A Q = \begin{bmatrix} q_1^T A q_1 & \beta^T \\ \alpha & C \end{bmatrix}, \quad \text{其中 } \alpha = Q_{n-1}^T A q_1 \in \mathbb{R}^{n-1},$$

则有

$$\|\alpha\|_2 = \|A q_1 - (q_1^T A q_1) q_1\|_2. \quad (7.53)$$

证明 因为对任意的 $x \in \mathbb{R}^n$ 都有 $\|Q^T x\|_2 = \|x\|_2$. 故

$$\begin{aligned} \|A q_1 - (q_1^T A q_1) q_1\|_2 &= \|[q_1, Q_{n-1}]^T [A q_1 - (q_1^T A q_1) q_1]\|_2 \\ &= \left\| \begin{bmatrix} q_1^T [A q_1 - (q_1^T A q_1) q_1] \\ Q_{n-1}^T [A q_1 - (q_1^T A q_1) q_1] \end{bmatrix} \right\|_2 = \left\| \begin{bmatrix} 0 \\ \alpha \end{bmatrix} \right\|_2 = \|\alpha\|_2. \end{aligned}$$

证毕. □

定理 7.13 设对称矩阵 $A \in \mathbb{R}^{n \times n}$ 满足

$$|\lambda_1| > |\lambda_2| \geq \cdots \geq |\lambda_n| > 0,$$

对应的规范正交化特征向量为 x_1, x_2, \dots, x_n . 如果单位坐标向量

$$e_1 = (1, 0, \dots, 0)^T = \sum_{i=1}^n \alpha_i x_i$$

中的 $\alpha_1 \neq 0$, 那么由算法 7.5 生成的矩阵序列 $\{A_k\}$ 具有收敛性质

$$\lim_{k \rightarrow \infty} A_k e_1 = \lambda_1 e_1. \quad (7.54)$$

证明 记 \tilde{Q}_k 的第 1 列为 $\tilde{q}_1^{(k)} = \tilde{Q}_k e_1$, \tilde{R}_k 的第 1 个对角元为 $\tilde{r}_{11}^{(k)}$, 则由式 (7.52) 可知

$$A^k e_1 = \tilde{Q}_k \tilde{R}_k e_1 = \tilde{Q}_k (\tilde{r}_{11}^{(k)} e_1) = \tilde{r}_{11}^{(k)} \tilde{q}_1^{(k)}.$$

注意到 $\|\tilde{q}_1^{(k)}\|_2 = 1$, $\tilde{r}_{11}^{(k)} \neq 0$ (因矩阵 A 非奇异). 从而 $|\tilde{r}_{11}^{(k)}| = \|A^k e_1\|_2$. 于是, 根据幂法的收敛性, 有

$$\lim_{k \rightarrow \infty} \tilde{q}_1^{(k)} = \lim_{k \rightarrow \infty} \frac{A^k e_1}{\tilde{r}_{11}^{(k)}} = z_1, \quad (7.55)$$

式中: z_1 为矩阵 A 的对应于 λ_1 的规范化特征向量 (可以相差一个常数因子 ± 1). 进而, 由式 (7.51), 可以把 A_{k+1} 写成

$$\begin{aligned} A_{k+1} &= \tilde{Q}_k^T A \tilde{Q}_k = [\tilde{q}_1^{(k)}, \tilde{Q}_{k-1}]^T A [\tilde{q}_1^{(k)}, \tilde{Q}_{k-1}] \\ &= \begin{bmatrix} a_{11}^{(k+1)} & * \\ \alpha^{(k+1)} & * \end{bmatrix}, \end{aligned}$$

式中: $a_{11}^{(k+1)} = (\tilde{q}_1^{(k)})^T A \tilde{q}_1^{(k)}$, $\alpha^{(k+1)} = (\tilde{Q}_{k-1})^T A \tilde{q}_1^{(k)}$. 根据式 (7.55), 得

$$\lim_{k \rightarrow \infty} a_{11}^{(k+1)} = \lim_{k \rightarrow \infty} (\tilde{q}_1^{(k)})^T A \tilde{q}_1^{(k)} = z_1^T A z_1 = \lambda_1.$$

故由引理 7.1, 得

$$\begin{aligned} \lim_{k \rightarrow \infty} \|\alpha^{(k+1)}\|_2 &= \lim_{k \rightarrow \infty} \|A \tilde{q}_1^{(k)} - [(\tilde{q}_1^{(k)})^T A \tilde{q}_1^{(k)}] \tilde{q}_1^{(k)}\|_2 \\ &= \|A z_1 - [z_1^T A z_1] z_1\|_2 = 0. \end{aligned}$$

因此, $\lim_{k \rightarrow \infty} \alpha^{(k+1)} = 0$. 于是有

$$\lim_{k \rightarrow \infty} A_{k+1} = \begin{bmatrix} \lambda_1 & * \\ 0 & * \end{bmatrix},$$

即式 (7.54) 成立. 证毕. □

作为定理 7.13 的推广, 有下面的收敛性结果.

定理 7.14 设 $A = X\Lambda X^{-1}$, 其中 $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$. 如果 ① $|\lambda_1| > |\lambda_2| > \dots > |\lambda_n| > 0$; ② X^{-1} 具有 LU 分解: $X^{-1} = LU$. 则 QR 方法基本收敛于上三角矩阵.

证明 由式 (7.51) 可知, 只需分析 \tilde{Q}_k 的极限情况. 而由式 (7.52), 只需分析 A^k 的极限情况. 注意到

$$A^k = X\Lambda^k X^{-1} = X\Lambda^k LU = X(\Lambda^k L\Lambda^{-k})\Lambda^k U.$$

令

$$\Lambda^k L\Lambda^{-k} = I + E_k,$$

则

$$A^k = X(I + E_k)\Lambda^k U.$$

由于 L 是单位下三角形, 故

$$(E_k)_{ij} = \begin{cases} 0, & i \leq j, \\ l_{ij}(\lambda_i/\lambda_j)^k, & i > j. \end{cases}$$

由假设条件 ① 知, $E_k \rightarrow O$ 且 $(E_k)_{ij}$ 的收敛速度是 $|\lambda_i/\lambda_j|$.

设 $X = QR$ 且 R 的对角元均为正数, 则有

$$\begin{aligned} A^k &= QR(I + E_k)\Lambda^k U \\ &= Q(I + RE_k R^{-1})R\Lambda^k U. \end{aligned}$$

因为 $E_k \rightarrow O$ ($k \rightarrow \infty$), 故当 k 充分大时, $I + RE_k R^{-1}$ 非奇异, 所以有唯一的 QR 分解 $I + RE_k R^{-1} = \hat{Q}_k \hat{R}_k$ (\hat{R}_k 的对角元为正), 而且当 $k \rightarrow \infty$ 时, $\hat{Q}_k \rightarrow I$, $\hat{R}_k \rightarrow I$. 此时, A^k 有如下分解

$$A^k = (Q\hat{Q}_k)(\hat{R}_k R\Lambda^k U).$$

妨碍上式成为 A^k 的 QR 分解的仅仅是上式右端第 2 个因子 (上三角矩阵) 的对角元可能非正. 为补救这一点, 可引入两个对角正交矩阵

$$D_1 = \text{diag}\left(\frac{\lambda_1}{|\lambda_1|}, \dots, \frac{\lambda_n}{|\lambda_n|}\right), \quad D_2 = \text{diag}\left(\frac{U_{11}}{|U_{11}|}, \dots, \frac{U_{nn}}{|U_{nn}|}\right),$$

式中: U_{ii} ($i = 1, 2, \dots, n$) 为矩阵 U 的对角元. 于是

$$A^k = ((Q\hat{Q}_k)(D_2 D_1^k))(D_1^{-k} D_2^{-1} \hat{R}_k R\Lambda^k U)$$

是 A^k 的唯一 QR 分解. 从而由式 (7.51), 有

$$A_{k+1} = (Q\hat{Q}_k D_2 D_1^k)^T A (Q\hat{Q}_k D_2 D_1^k).$$

将 $A = X\Lambda X^{-1} = QR\Lambda R^{-1}Q^{-1}$ 代入上式, 得

$$A_{k+1} = (D_2 D_1^k)^T (\hat{Q}_k^{-1} R\Lambda R^{-1} \hat{Q}_k) (D_2 D_1^k).$$

因为 $\hat{Q}_k^{-1} R\Lambda R^{-1} \hat{Q}_k \rightarrow R\Lambda R^{-1} \equiv \bar{R}$ (上三角矩阵), 所以 A_k 的对角线以下元素收敛于 0. 因为 D_1^k 可能不收敛, 故 A_k 基本收敛于 \bar{R} . 证毕. \square

7.4.4 带原点位移的 QR 方法

定理 7.14 表明 \mathbf{A}_k 的对角元 $a_{ii}^{(k)} \rightarrow \lambda_i (k \rightarrow \infty)$. 从证明过程中可以看出 \mathbf{A}_k 的下三角部分的元素趋于 0 的速度由 $k \rightarrow \infty$ 时, $\hat{\mathbf{Q}}_k \rightarrow \mathbf{I}$ 和 $\hat{\mathbf{R}}_k \rightarrow \mathbf{I}$ 的速度, 亦即由 $\mathbf{E}_k \rightarrow \mathbf{O}$ 的速度所决定. 而从矩阵 \mathbf{E}_k 的构成可以看出, \mathbf{E}_k 的第 i 行元素趋于 0 的速度由 $|\lambda_i/\lambda_{i-1}|$ 决定, 第 i 列元素趋于 0 的速度由 $|\lambda_{i+1}/\lambda_i|$ 决定. 可以证明, \mathbf{A}_k 的下三角部分的元素趋于 0 的情况也是这样. 所以 $a_{ii}^{(k)} \rightarrow \lambda_i$ 的速度由 \mathbf{A}_k 的第 i 行和第 i 列的下三角部分的元素趋于 0 的速度确定, 这个速度即为 $O(\rho_i)$, 其中

$$\rho_i = \max \left\{ \frac{|\lambda_i|}{|\lambda_{i-1}|}, \frac{|\lambda_{i+1}|}{|\lambda_i|}, \lambda_0 = +\infty, \lambda_{n+1} = 0 \right\}.$$

在实际计算中, 线性收敛速度是不令人满意的, 特别是当 ρ_i 不算很小时, 收敛将是十分缓慢的. 为此, 将使用原点位移的方法进行加速. 现在 $\rho_n = |\lambda_n/\lambda_{n-1}|$, 如果将算法用于矩阵 $\mathbf{A} - s\mathbf{I}$, 则 $a_{nn}^{(k)}$ 将以商 $|\lambda_n - s|/|\lambda_{n-1} - s|$ 线性收敛于 $\lambda_n - s$. 当 s 是 λ_n 的一个较好的近似时, 收敛是很快的. 基于这个想法, 可构造原点位移 QR 方法如下.

算法 7.6 (原点位移 QR 方法)

步 1, 输入上 Hessenberg 矩阵 $\mathbf{A} \in \mathbb{R}^{n \times n}$.

步 2, 记 $\mathbf{A}_1 := \mathbf{A}$. 对于 $k = 1, 2, \dots$,

$$(1) \quad \mathbf{A}_k - s_k \mathbf{I} = \mathbf{Q}_k \mathbf{R}_k \quad (\text{QR 分解}).$$

$$(2) \quad \mathbf{A}_{k+1} = \mathbf{Q}_k^T \mathbf{A}_k \mathbf{Q}_k = \mathbf{R}_k \mathbf{Q}_k + s_k \mathbf{I} \quad (\text{正交相似变换}).$$

其中, 由 \mathbf{A}_k 到 \mathbf{A}_{k+1} 的变换称为原点位移的 QR 变换.

与算法 7.5 相类似, 由此生成的矩阵序列 $\{\mathbf{A}_k\}$ 和 $\{\mathbf{Q}_k\}$ 都是上 Hessenberg 形, 并且 \mathbf{A}_k 与 \mathbf{A} 相似. 现在的问题是如何选择 s_k 使得收敛速度加快. 下面讨论一类特殊情形.

假设 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 满足定理 7.13 的条件, 并设由基本 QR 方法生成的 \mathbf{A}_k 右下角 2 阶子矩阵

$$\mathbf{J}_k = \begin{bmatrix} a_{n-1,n-1}^{(k)} & a_{n-1,n}^{(k)} \\ a_{n,n-1}^{(k)} & a_{n,n}^{(k)} \end{bmatrix} \quad (7.56)$$

的特征值为 $\lambda_{n-1}^{(k)}$ 和 $\lambda_n^{(k)}$, 则有

$$\lim_{k \rightarrow \infty} a_{n,n}^{(k)} = \lambda_n, \quad \lim_{k \rightarrow \infty} \lambda_{n-1}^{(k)} = \lambda_{n-1}, \quad \lim_{k \rightarrow \infty} \lambda_n^{(k)} = \lambda_n. \quad (7.57)$$

下面考虑位移量 s_k 的取法. 因为矩阵 $\mathbf{A}_k - s_k \mathbf{I}$ 的特征值 $\mu_k = \lambda_k - s_k$, 因此, 如果取 $s_k \approx \lambda_n$, 那么在定理 7.13 的条件下, 有

$$\left| \frac{\mu_{k+1}}{\mu_k} \right| = \left| \frac{\lambda_{k+1} - \lambda_n}{\lambda_k - \lambda_n} \right| < \left| \frac{\lambda_{k+1}}{\lambda_k} \right|, \quad k = 1, 2, \dots, n-1.$$

这表明算法 7.6 的收敛速度比算法 7.5 快. 因此, 由子矩阵 (7.56) 的收敛性质 (7.57), 位移量 s_k 有下列两种取法:

$$(1) s_k = a_{nn}^{(k)}.$$

(2) 当 $\lambda_{n-1}^{(k)}$ 和 $\lambda_n^{(k)}$ 为实数时, 取 s_k 为其中与 $a_{nn}^{(k)}$ 最接近的一个.

这两种位移策略特别适用于对称矩阵, 因为此时子矩阵 J_k 对称, 从而它的两个特征值都是实数. 可以证明, 采用这种位移策略的算法 7.6, A_k 基本收敛于上三角矩阵, 收敛是二阶的, 并且 $a_{n,n-1}^{(k)}$ 最先趋于零, 从而首先得到绝对值最小的特征值 λ_n .

7.4.5 双重步位移隐式 QR 方法

对于一般的实方阵 $A \in \mathbb{R}^{n \times n}$, 前面讨论的算法 7.5 和算法 7.6 都不太好. 首先, 这两种方法都是显式方法, 即每次迭代都需要明显地作矩阵的 QR 分解, 并用所得到的正交矩阵作相似变换 (实际上是作三角矩阵与正交矩阵的乘积), 计算量和存储量都很大. 其次, 基本 QR 算法 7.5 若收敛则是线性的, 而原点位移 QR 算法 7.6 只有当特征值是实数时才有可能改善收敛性. 这是可以理解的, 因为这两种算法都是在实数域上进行运算, 因此当实矩阵 A 有复特征值时, 即使收敛也是很缓慢的.

理论分析和实际计算的实验表明: QR 迭代产生的矩阵序列右下角最先显露 A 的特征值. 在原点位移 QR 方法中正是利用这一特点来选取位移参数 s_k 的. 如果显露的是 A 的实特征值, 即 $a_{nn}^{(k)}$ 是 A 的较好的近似特征值时, 就可以简单地选取 $s_k = a_{nn}^{(k)}$. 然而, 如果显露的是 A 的复共轭特征值时, 即式 (7.56) 中的矩阵 J_k 的特征值是一对互相共轭的复数 μ_1 和 μ_2 , 且与 A 的特征值比较接近时, 就应该选择 J_k 某一特征值 μ_i 作为位移参数. 但这样一来, 就引进了复运算, 而这是所不希望的.

实数运算的优点是, 计算简单, 工作量小. 为了用实数运算来求一般实方阵的全部特征值, 特别是复特征值, 可以引进双重步位移隐式 QR 方法. 它的基本思想是: 首先把原点位移推广到复数域 (当然包括实数), 并且每次迭代作两次原点位移的 QR 变换, 因此称为双重步 QR 方法, 这时的运算都可以是复数; 然后把这两次 QR 变换合在一起, 转化成实数运算, 并且构成隐式方法, 即不明显地进行 QR 变换. 它的优点是, 迭代过程都是实数运算, 原点位移加速了收敛性, 隐式方法可减少计算工作量并节省存储空间.

1. 双重步位移隐式 QR 变换

双重步位移隐式 QR 方法的核心是双重步位移隐式 QR 变换. 设矩阵 $A \in \mathbb{R}^{n \times n}$ 是不可约的上 Hessenberg 矩阵, 它的右下角的 2 阶子矩阵

$$\begin{bmatrix} a_{n-1,n-1} & a_{n-1,n} \\ a_{n,n-1} & a_{n,n} \end{bmatrix}$$

的特征值为 μ_1 和 μ_2 , 它们可能都是实数, 也可能是一对共轭复数. 记

$$s = a_{n-1,n-1} + a_{n,n}, \quad t = a_{n-1,n-1}a_{n,n} - a_{n,n-1}a_{n-1,n},$$

s 和 t 都是实数, 则有

$$\mu_1 + \mu_2 = s, \quad \mu_1 \mu_2 = t. \quad (7.58)$$

现在取 μ_1 和 μ_2 作为位移量, 作两次原点位移的 QR 变换:

$$\begin{cases} A - \mu_1 I = Q_1 R_1, & B = Q_1^H A Q_1 = R_1 Q_1 + \mu_1 I, \\ B - \mu_2 I = Q_2 R_2, & C = Q_2^H B Q_2 = R_2 Q_2 + \mu_2 I, \end{cases} \quad (7.59)$$

式中: Q_1, Q_2 为酉矩阵, 也是上 Hessenberg 矩阵; B 和 C 都是上 Hessenberg 矩阵; R_1 和 R_2 都是上三角矩阵. 双重步位移的 QR 变换 (7.59) 具有如下性质: 若记

$$H = (A - \mu_2 I)(A - \mu_1 I), \quad Q = Q_1 Q_2, \quad R = R_2 R_1, \quad (7.60)$$

式中: Q 为酉矩阵; R 为上三角矩阵. 则有

$$H = QR, \quad C = Q^H A Q. \quad (7.61)$$

事实上, 有

$$\begin{aligned} H &= (A - \mu_2 I)(A - \mu_1 I) = (A - \mu_2 I)Q_1 R_1 \\ &= Q_1(Q_1^H A Q_1 - \mu_2 I)R_1 = Q_1(B - \mu_2 I)R_1 \\ &= Q_1(Q_2 R_2)R_1 = QR, \\ C &= Q_2^H B Q_2 = Q_2^H(Q_1^H A Q_1)Q_2 = Q^H A Q. \end{aligned}$$

进一步, 注意到

$$\begin{aligned} H &= A^2 - (\mu_1 + \mu_2)A + \mu_1 \mu_2 I = A^2 - sA + tI \\ &= \begin{bmatrix} h_{11} & \times & \times & \times & \cdots & \times \\ h_{21} & \times & \times & \times & \cdots & \times \\ h_{31} & \times & \times & \times & \cdots & \times \\ & \times & \times & \times & \cdots & \times \\ & & \ddots & \ddots & \ddots & \vdots \\ & & & \times & \times & \times \end{bmatrix} \end{aligned} \quad (7.62)$$

是实矩阵, 其中

$$\begin{cases} h_{11} = a_{11}^2 + a_{12}a_{21} - sa_{11} + t, \\ h_{21} = a_{21}(a_{11} + a_{22} - s), \\ h_{31} = a_{21}a_{32}. \end{cases} \quad (7.63)$$

所以式 (7.61) 的第 1 式 $H = QR$ 是实的 QR 分解, 即其中的酉矩阵 Q 必定是正交矩阵, R 必是实上三角矩阵. 从而式 (7.61) 的第 2 式为正交相似变换 $C = Q^T A Q$.

综上所述, 为了确保计算得到的 C 仍为实矩阵, 根据式 (7.61), 自然考虑按如下的步骤来计算 C :

(1) 计算 $H = A^2 - sA + tI$.

(2) 计算 H 的 QR 分解: $H = QR$.

(3) 计算 C 的正交相似变换: $C = Q^T A Q$.

然而, 如此计算的第 1 步形成 H 的运算量就为 $O(n^3)$, 这使得前面为减少每次迭代所需运算量所做的努力付之东流. 幸运的是, 定理 7.12 表明, 不论采取什么方法计算正交矩阵 \tilde{Q} 使得 $\tilde{Q}^T A \tilde{Q} = \tilde{C}$ 成为上 Hessenberg 矩阵, 只要保证 \tilde{Q} 的第 1 列与 Q 的第 1 列一样, 则 \tilde{C} 就与 C 在本质上是一样的 (所有元素的绝对值相等). 因此, 可以有很大的自由度去寻求更有效的方法来实现 A 到 C 的变换.

首先, 从式 (7.61) 的第 1 式 $H = QR$ 知, Q 的第 1 列与 H 的第 1 列共线, 即 Qe_1 由 He_1 单位化得到. 而由式 (7.62) 容易算出

$$He_1 = (h_{11}, h_{21}, h_{31}, 0, \dots, 0)^T,$$

式中: h_{11}, h_{21}, h_{31} 由式 (7.63) 得出.

其次, 如果 Householder 变换 P_0 将 He_1 变为 αe_1 , 即 $P_0(He_1) = \alpha e_1$, 其中 $\alpha \in \mathbb{R}$, 则易知, P_0 的第 1 列就与 He_1 共线, 从而 $P_0 e_1 = Qe_1$. 而由 Householder 变换的性质, P_0 可按如下方式确定:

$$P_0 = \text{diag}(\tilde{P}_0, I_{n-3}),$$

式中:

$$\tilde{P}_0 = I_3 - \beta vv^T, \quad \beta = 2/v^T v,$$

$$v = (h_{11} + \alpha \text{sign}(h_{11}), h_{21}, h_{31})^T, \quad \alpha = \sqrt{h_{11}^2 + h_{21}^2 + h_{31}^2}.$$

现令

$$D_1 = P_0 A P_0,$$

那么只要找到第 1 列为 e_1 的正交矩阵 \tilde{Q} 使 $\tilde{Q}^T D_1 \tilde{Q} = \tilde{H}$ 为上 Hessenberg 矩阵, 那么 \tilde{H} 就是希望得到的矩阵 C . 这只需确定 $n-1$ 个 Householder 变换 P_1, \dots, P_{n-1} , 使

$$(P_{n-1} \cdots P_1) D_1 (P_1 \cdots P_{n-1}) = \tilde{H}$$

为上 Hessenberg 矩阵, 即有 $\tilde{Q} = P_1 \cdots P_{n-1}$ 的第 1 列为 e_1 . 而且由于 D_1 所具有的特殊性质, 实现这一约化过程所需的运算量仅为 $O(n^2)$.

事实上, 由于用 P_0 将 A 相似变换为 D_1 只改变了 A 的前三行和前三列, 故 D_1 具有如下形状, 即

$$D_1 = P_0 A P_0 = \begin{bmatrix} \times & \times & \times & \times & \cdots & \times & \times \\ \times & \times & \times & \times & \cdots & \times & \times \\ \oplus & \times & \times & \times & \cdots & \times & \times \\ \oplus & \oplus & \times & \times & \cdots & \times & \times \\ & & & \times & \cdots & \vdots & \vdots \\ & & & & \ddots & \vdots & \vdots \\ & & & & & \times & \times \end{bmatrix},$$

仅比上 Hessenberg 形多 3 个非零元“ \oplus ”. 由 D_1 的这种特殊性, 易知用来约化 D_1 为上 Hessenberg 形的第一个 Householder 变换 P_1 具有如下形状, 即

$$P_1 = \text{diag}(1, \tilde{P}_1, I_{n-4}),$$

式中: \tilde{P}_1 为 3 阶 Householder 变换, 而且 $P_1 D_1 P_1$ 具有如下形状, 即

$$D_2 = P_1 D_1 P_1 = P_1 P_0 A P_0 P_1 = \begin{bmatrix} \times & \times & \times & \times & \cdots & \times & \times \\ \times & \times & \times & \times & \cdots & \times & \times \\ 0 & \times & \times & \times & \cdots & \times & \times \\ 0 & \oplus & \times & \times & \cdots & \times & \times \\ & \oplus & \oplus & \times & \cdots & \vdots & \vdots \\ & & & & \ddots & \vdots & \vdots \\ & & & & & \times & \times \end{bmatrix},$$

如此递推地进行下去, 不难发现, 第 k 次约化所用的 Householder 变换 P_k 具有如下形状, 即

$$P_k = \text{diag}(I_k, \tilde{P}_k, I_{n-k-3}),$$

式中: \tilde{P}_k 为 3 阶 Householder 变换, $k = 2, \cdots, n-3$, 而且 $D_{n-2} := P_{n-3} D_{n-2} P_{n-3}$ 具有如下形状

$$\begin{aligned} D_{n-2} &= P_{n-3} D_{n-3} P_{n-3} \\ &= P_{n-3} P_{n-4} \cdots P_1 D_1 P_1 \cdots P_{n-4} P_{n-3} \\ &= \begin{bmatrix} \times & \times & \times & \cdots & \times & \times \\ \times & \times & \times & \cdots & \times & \times \\ & \times & \times & \cdots & \times & \times \\ & & & \ddots & \vdots & \vdots \\ & & & & \times & \times \\ & & & & \oplus & \times & \times \end{bmatrix}. \end{aligned}$$

因此, 最后一次约化所用的 Householder 变换 P_{n-2} 具有如下形状, 即

$$P_{n-1} = \text{diag}(I_{n-2}, \tilde{P}_{n-2}),$$

式中: \tilde{P}_{n-2} 为 2 阶 Householder 变换. 而且 $D_{n-1} := P_{n-2} D_{n-2} P_{n-2}$ 具有如下形状, 即

$$\begin{aligned} D_{n-1} &= P_{n-2} D_{n-2} P_{n-2} \\ &= P_{n-2} \cdots P_1 P_0 A P_0 P_1 \cdots P_{n-2} \end{aligned}$$

$$= \begin{bmatrix} \times & \times & \times & \cdots & \times & \times \\ \times & \times & \times & \cdots & \times & \times \\ & \times & \times & \cdots & \times & \times \\ & & \ddots & \ddots & \vdots & \vdots \\ & & & \times & \times & \times \\ & & & & \times & \times \end{bmatrix} := D. \quad (7.64)$$

这样, 就找到了一种由 A 到 C 的变换方法, 它既避免了复运算的出现, 又减少了运算量. 当然, 这一变换过程对 J_k 的两个特征值都是实数的情形也是可行的. 因此, 不必在选取位移参数时区别显露的是实特征值还是复共轭特征值的情形, 而只需取作 J_k 的两个特征值即可.

定理 7.15 由式 (7.64) 得到的矩阵 D “基本收敛”于变换 (7.61) 中的矩阵 C .

证明 若记正交矩阵

$$U = P_0 P_1 \cdots P_{n-2},$$

则式 (7.64) 即为 $D = U^T A U$. 易知, U 的第 1 列与 P_0 的第 1 列相同, 即

$$U e_1 = P_0 e_1.$$

另外, 用 Householder 变换对矩阵 H 作 QR 分解的过程是

$$\tilde{P}_{n-1} \cdots \tilde{P}_2 \tilde{P}_1 H = R,$$

式中: $\tilde{P}_1 = P_0$. 对于 $i = 2, \cdots, n-2$, 有

$$\tilde{P}_i = \begin{bmatrix} I_{i-1} & & \\ & \tilde{V}_i & \\ & & I_{n-i-2} \end{bmatrix},$$

式中: $\tilde{V}_i \in \mathbb{R}^{3 \times 3}$ 为 Householder 矩阵.

当 $i = n-1$ 时, 有

$$\tilde{P}_{n-1} = \begin{bmatrix} I_{n-2} & \\ & \tilde{V}_{n-1} \end{bmatrix},$$

式中: $\tilde{V}_{n-1} \in \mathbb{R}^{2 \times 2}$ 为 Householder 矩阵. 若记正交矩阵

$$Q = \tilde{P}_1 \tilde{P}_2 \cdots \tilde{P}_{n-1} = P_0 \tilde{P}_2 \cdots \tilde{P}_{n-1},$$

则由 $H = QR$. 注意到 Q 的第 1 列

$$Q e_1 = P_0 e_1 = U e_1,$$

根据定理 7.12, 在化上 Hessenberg 形 $C = Q^T A Q$ 的变换中, 矩阵 Q 和 C 本质上由 Q 的第 1 列唯一确定 (相差一个 1 或 -1 因子的意义下). 因此, 定理的结论成立. 证毕. \square

综上所述, 可得双重步位移的 QR 变换的迭代过程.

算法 7.7 (双重步位移隐式 QR 变换)

步 1, 输入不可约上 Hessenberg 矩阵 $A = (a_{ij}) \in \mathbb{R}^{n \times n}$.

步 2, $k := 0; m := n - 1; s := a_{mm} + a_{nn};$

$$t := a_{mm}a_{nn} - a_{mn}a_{nm};$$

$$x := a_{11}^2 + a_{12}a_{21} - sa_{11} + t;$$

$$y := a_{21}(a_{11} + a_{22} - s); z := a_{21}a_{32}.$$

步 3, 若 $k = n - 2$, 则转步 5; 否则, 确定 Householder 矩阵 $\tilde{P}_k \in \mathbb{R}^{3 \times 3}$ 使

$$\tilde{P}_k \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} * \\ 0 \\ 0 \end{bmatrix},$$

置 $A := P_k A P_k$, 其中 $P_k = \text{diag}(I_k, \tilde{P}_k, I_{n-k-3})$.

步 4, 更新:

$$x := a_{k+2,k+1}, \quad y := a_{k+3,k+1}, \quad z := \begin{cases} a_{k+4,k+1}, & k < n - 3 \\ 0, & k = n - 3. \end{cases}$$

置 $k := k + 1$, 转步 3.

步 5 确定 Householder 矩阵 $\tilde{P}_{n-2} \in \mathbb{R}^{2 \times 2}$ 使

$$\tilde{P}_{n-2} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} * \\ 0 \end{bmatrix},$$

置 $A := P_{n-2} A P_{n-2}$, 其中 $P_k = \text{diag}(I_{n-2}, \tilde{P}_{n-2})$. 迭代结束.

算法 7.7 的运算量是 $6n^2$. 如果需要把正交变换累积起来, 还需再增加运算量 $6n^2$.

双重步位移隐式 QR 变换的 MATLAB 程序如下:

```
function A=ddi_qrtran(A,n)
%双重位移隐式QR变换
%输入:A为不可约n阶上Hessenberg矩阵,其最后2x2阶主子阵有特征值a和b
%输出:A=Q^T A Q,这里Q=P1...P_{n-2}是一系列Householder矩阵的乘积
%    且Q^T(A-aI)(A-bI)是上三角形矩阵
I3=eye(3); I2=eye(2); s=A(n-1,n-1)+A(n,n);
```

```

t=A(n-1,n-1)*A(n,n)-A(n-1,n)*A(n,n-1);
x=A(1,1)^2+A(1,2)*A(2,1)-s*A(1,1)+t;
y=A(2,1)*(A(1,1)+A(2,2)-s); z=A(2,1)*A(3,2);
for k=0:n-3
    [v,beta]=mhouse([x,y,z]'); q=max(1,k);
    A(k+1:k+3,q:n)=(I3-beta*v*v')*A(k+1:k+3,q:n);
    r=min(k+4,n);
    A(1:r,k+1:k+3)=A(1:r,k+1:k+3)*(I3-beta*v*v');
    x=A(k+2,k+1); y=A(k+3,k+1);
    if (k<n-3), z=A(k+4,k+1); end
end
[v,beta]=mhouse([x,y]');
A(n-1:n,n-2:n)=(I2-beta*v*v')*A(n-1:n,n-2:n);
A(1:n,n-1:n)=A(1:n,n-1:n)*(I2-beta*v*v');

```

2. 双重步位移隐式 QR 方法

前面的讨论已经解决了用 QR 方法求一个给定实矩阵的实 Schur 分解的几个关键性问题. 然而, 作为一种实用的算法, 还需给出一种有效的判定准则, 来判定迭代过程中所产生的上 Hessenberg 矩阵的次对角元素何时可以忽略不计. 一种简单而实用的准则是: 当

$$|a_{i+1,i}| \leq (|a_{ii}| + |a_{i+1,i+1}|)\varepsilon \quad (7.65)$$

时, 就将 $a_{i+1,i}$ 看作是 0.

将算法 7.4 和算法 7.7 与收敛准则 (7.65) 结合起来, 就得到了双重步位移隐式 QR 方法. 这一算法是计算一给定的 n 阶实矩阵 A 实 Schur 分解: $Q^T A Q = T$, 其中 Q 为正交矩阵, T 为拟上三角矩阵, 即对角块为 1×1 或 2×2 方阵的块上三角矩阵. 由于其中的 QR 变换不明显, 故称为隐式方法.

算法 7.8 (双重步位移隐式 QR 方法)

- 步 1, 输入矩阵 $A = (a_{ij}) \in \mathbb{R}^{n \times n}$.
- 步 2, 上 Hessenberg 化. 用算法 7.4 计算 A 的上 Hessenberg 分解 $A := Q^T A Q$.
- 步 3, 若满足收敛性准则, 停算; 否则返回步 2 继续进行双重步位移隐式 QR 迭代.

双重步位移隐式 QR 方法的 MATLAB 程序如下:

```

function [iter,D]=ddiqr_eig(A,tol)
%用双重步位移隐式QR方法求实方阵的全部特征值
%输入:A为n阶上Hessenberg形实方阵,tol为控制精度(默认是1e-5)
%输出:iter为迭代次数,D为A的全部特征值

```

```

if nargin<2, tol=1e-5; end
n=size(A,1);
D=zeros(n,1); i=n; m=n; iter=0; %初始化
[A]=hessenb(A); %化矩阵A为Hessenberg矩阵
while (m>0)
    %用双重位移隐式QR方法进行迭代
    if m<=2
        la=eig(A(1:m,1:m));
        D(1:m)=la'; break;
    end
    iter=iter+1;
    A=ddiqr_tran(A,m); %对上Hessenberg 矩阵作QR分解,并作正交相似变换
    for k=m-1:-1:1 %下面的程序段是判断是否终止
        if abs(A(k+1,k))<tol
            if m-k<=2
                la=eig(A(k+1:m,k+1:m));
                j=i-m+k+1; D(j:i)=la';
                i=j-1; m=k; break;
            end
        end
    end
end
end
end

```

例 7.8 用双重步位移隐式 QR 方法求例 7.7 中矩阵 A 的全部特征值.

解 编写 M 文件 ex78.m, 在 MATLAB 命令窗口输入 ex78 得

```

>> ex78

```

算 法	迭代次数	CPU时间
基本QR方法	622	0.0702
双位移QR方法	8	0.0201

基本QR方法特征值	双位移QR方法特征值
18.4123 + 0.0000i	18.4123 + 0.0000i
11.1805 + 0.0000i	11.1805 + 0.0000i
1.7099 - 4.2522i	1.7099 - 4.2522i
1.7099 + 4.2522i	1.7099 + 4.2522i
4.4983 + 0.0000i	-2.2327 + 0.0000i
-2.2327 + 0.0000i	4.4983 + 0.0000i
-0.2783 + 0.0000i	-0.2783 + 0.0000i

7.4.6 特征向量的计算方法

本节讨论在用 (双重步位移隐式) QR 方法求得给定矩阵的特征值之后, 如何求其

对应的特征向量. 设 $A \in \mathbb{R}^{n \times n}$, 并假定已用 QR 方法求得 A 的特征值 λ 的一个近似 $\tilde{\lambda}$. 现在讨论如何求对应于 λ 的特征向量.

目前, 解决这一问题最好的方法是带原点位移的反幂法 (也称为反迭代法, 逆迭代法), 其基本迭代格式如下:

$$(A - \alpha I)y^{(k)} = x^{(k-1)}, \quad (7.66a)$$

$$x^{(k)} = y^{(k)} / \|y^{(k)}\|_2, \quad k = 1, 2, \dots, \quad (7.66b)$$

式中: α 为选定的位移参数; $x^{(0)}$ 为给定的初始向量.

从式 (7.66a) 可以看出, 每迭代一次就需要解一个线性方程组, 这要比幂法的运算量大得多. 但由于方程组的系数矩阵不随 k 变化, 故可事先对它进行列主元 LU 分解, 然后每次迭代只需解两个三角形方程组即可. 顺便指出, 式 (7.66b) 只是为了防止迭代“溢出”而做的归一化处理, 在实际计算时可以用 $\|\cdot\|_\infty$ 进行归一化.

现假定 A 是非亏损的, 即存在 $X = [\xi_1, \xi_2, \dots, \xi_n] \in \mathbb{C}^{n \times n}$ 非奇异, 使得

$$X^{-1}AX = A \equiv \text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_n\}. \quad (7.67)$$

而且不失一般性, 还可以假设 $\|\xi_i\|_2 = 1$ ($i = 1, 2, \dots, n$). 现将初始向量 $x^{(0)}$ 按 $\xi_1, \xi_2, \dots, \xi_n$ 展开

$$x^{(0)} = \sum_{i=1}^n \beta_i \xi_i. \quad (7.68)$$

再假定 α 与 A 的特征值 λ_s 最靠近, 并有

$$0 < |\alpha - \lambda_s| < |\alpha - \lambda_i|, \quad i \neq s, \quad (7.69a)$$

$$\beta_s \neq 0. \quad (7.69b)$$

由式 (7.66a)、式 (7.67) 和式 (7.68), 得

$$\begin{aligned} y^{(k)} &= \theta_k (A - \alpha I)^{-k} x^{(0)} \\ &= \theta_k \sum_{i=1}^n \beta_i (A - \alpha I)^{-k} \xi_i = \theta_k \sum_{i=1}^n \beta_i (\lambda_i - \alpha)^{-k} \xi_i \\ &= \theta_k \beta_s (\lambda_s - \alpha)^{-k} \left[\xi_s + \sum_{i \neq s}^n \left(\frac{\lambda_s - \alpha}{\lambda_i - \alpha} \right)^k \frac{\beta_i}{\beta_s} \xi_i \right] \\ &:= \theta_k \beta_s (\lambda_s - \alpha)^{-k} (\xi_s + u_k), \end{aligned} \quad (7.70)$$

式中: θ_k 为正数, 而 $u_k \rightarrow 0$ ($k \rightarrow \infty$), 其收敛速度依赖于

$$\frac{|\lambda_s - \alpha|}{\min_{i \neq s} |\alpha - \lambda_i|}$$

的大小. 将式 (7.70) 代入式 (7.66b), 得

$$x^{(k)} = \frac{y^{(k)}}{\|y^{(k)}\|_2} = \frac{\eta_k}{\|\xi_s + u_k\|_2} (\xi_s + u_k),$$

式中:

$$\eta_k = \frac{\theta_k \beta_s (\lambda_s - \alpha)^{-k}}{|\theta_k \beta_s (\lambda_s - \alpha)^{-k}|} = \frac{\beta_s}{|\beta_s|} \left(\frac{\lambda_s - \alpha}{|\lambda_s - \alpha|} \right)^{-k}$$

为满足 $|\eta_k| = 1$ 的复数, 从而有

$$\begin{aligned} \text{dist}(\mathbf{x}^{(k)}, \boldsymbol{\xi}_s) &= \|\mathcal{P}_{\mathbf{x}^{(k)}} - \mathcal{P}_{\boldsymbol{\xi}_s}\|_2 = \|\mathbf{x}^{(k)}(\mathbf{x}^{(k)})^H - \boldsymbol{\xi}_s \boldsymbol{\xi}_s^H\|_2 \\ &= \left\| \frac{(\boldsymbol{\xi}_s + \mathbf{u}_k)(\boldsymbol{\xi}_s + \mathbf{u}_k)^H}{(\boldsymbol{\xi}_s + \mathbf{u}_k)^H(\boldsymbol{\xi}_s + \mathbf{u}_k)} - \boldsymbol{\xi}_s \boldsymbol{\xi}_s^H \right\|_2 \\ &\rightarrow 0, \quad k \rightarrow \infty, \end{aligned}$$

收敛速度依赖于 $|\lambda_s - \alpha| / \min_{i \neq s} |\alpha - \lambda_i|$ 的大小. 换言之, 即 $\mathbf{x}^{(k)}$ 将按方向收敛于 \mathbf{A} 的特征向量. α 与 λ_s 越靠近, 收敛速度越快.

由此可见, 从收敛速度的角度来考虑, 用式 (7.66) 迭代时, 自然是 α 取得越靠近 \mathbf{A} 的某个特征值越好. 但是, 当 α 与 \mathbf{A} 的某个特征值很靠近时, $\mathbf{A} - \alpha \mathbf{I}$ 就很接近于一个奇异矩阵, 每一步迭代就需要解一个非常病态的线性方程组. 幸运的是理论分析以及大量的计算实践表明: $\mathbf{A} - \alpha \mathbf{I}$ 的病态性并不影响其收敛速度, 而且当 α 很靠近 \mathbf{A} 的某个特征值时, 常常只需要一次迭代就可以得到相当好的近似特征向量. 为此, 下面进行简要的理论分析.

设 λ 是 \mathbf{A} 的特征值, $\mathbf{y} \in \mathbb{C}^n$ 满足 $\|\mathbf{y}\|_2 = 1$. 定义

$$\mathbf{r} = (\mathbf{A} - \lambda \mathbf{I})\mathbf{y} \quad (7.71)$$

为向量 \mathbf{y} 的残差向量. 由式 (7.71), 得

$$(\mathbf{A} - \mathbf{r}\mathbf{y}^H)\mathbf{y} = \lambda\mathbf{y},$$

即 \mathbf{y} 是矩阵 $\mathbf{A} - \mathbf{r}\mathbf{y}^H$ 对应于 λ 的特征向量. 如果 $\|\mathbf{r}\|_2$ 很小, 则 $\|\mathbf{r}\mathbf{y}^H\|_2$ 也会很小. 从而, 当 $\|\mathbf{r}\|_2$ 很小时, \mathbf{y} 就是 \mathbf{A} 的对应于 λ 的一个很好的近似特征向量. 即可用 \mathbf{y} 的残差向量大小来衡量 \mathbf{y} 可否作为对应于 λ 的近似特征向量.

进一步, 假定

$$|\alpha - \lambda| = \min_{\tilde{\lambda} \in \lambda(\mathbf{A})} |\alpha - \tilde{\lambda}| \leq \varepsilon_1, \quad (7.72)$$

式中: ε_1 为很小的正数, 通常 $\varepsilon_1 = O(\epsilon)$; ϵ 为机器精度. 再假定对给定的初始向量 $\mathbf{x}^{(0)}$, 用列主元 Gauss 消去法求解方程组 (7.66a) 得到向量 $\mathbf{y}^{(1)}$ 的计算值 $\tilde{\mathbf{y}}^{(1)}$. 则由 Gauss 消去法的舍入误差分析结果可知

$$(\mathbf{A} - \alpha \mathbf{I} + \mathbf{E})\tilde{\mathbf{y}}^{(1)} = \mathbf{x}^{(0)}, \quad (7.73)$$

式中: $\|\mathbf{E}\|_2 \leq \varepsilon_2$ (通常 $\varepsilon_2 = O(\epsilon)$). 这样, 由式 (7.66b) 计算得

$$\mathbf{x}^{(1)} = \tilde{\mathbf{y}}^{(1)} / \|\tilde{\mathbf{y}}^{(1)}\|_2, \quad (7.74)$$

这里忽略了计算 $\mathbf{x}^{(1)}$ 所产生的误差.

利用式 (7.73), 可得向量 $\mathbf{x}^{(1)}$ 的残差向量为

$$\mathbf{r} = (\mathbf{A} - \lambda \mathbf{I})\mathbf{x}^{(1)} = (\alpha - \lambda)\mathbf{x}^{(1)} - \mathbf{E}\mathbf{x}^{(1)} + \frac{\mathbf{x}^{(0)}}{\|\tilde{\mathbf{y}}^{(1)}\|_2}.$$

于是有

$$\|\mathbf{r}\|_2 \leq \varepsilon_1 + \varepsilon_2 + \|\tilde{\mathbf{y}}^{(1)}\|_2^{-1}, \quad (7.75)$$

这里假定 $\|\mathbf{x}^{(0)}\|_2 = 1$.

由此可见, 如果计算得到的 $\tilde{\mathbf{y}}^{(1)}$ 具有很大的范数, 则由式 (7.66) 迭代一次所得的向量就有很小的残差向量, 从而在特征值问题不是十分病态的条件下, 就得到了很好的近似特征向量. 因此, 只需说明在式 (7.72) 成立的前提下确有 $\tilde{\mathbf{y}}^{(1)}$ 的范数很大.

事实上, 设 $\mathbf{A} - \alpha \mathbf{I} + \mathbf{E}$ 的奇异值分解为

$$\mathbf{A} - \alpha \mathbf{I} + \mathbf{E} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^H, \quad (7.76)$$

式中: $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n]$, $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n]$; $\mathbf{\Sigma} = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_n)$; $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n > 0$. 由 $\lambda - \alpha$ 是矩阵 $\mathbf{A} - \alpha \mathbf{I}$ 的特征值, 得

$$\sigma_n(\mathbf{A} - \alpha \mathbf{I}) \leq |\lambda - \alpha| \leq \varepsilon_1,$$

式中: $\sigma_n(\mathbf{A} - \alpha \mathbf{I})$ 为 $\mathbf{A} - \alpha \mathbf{I}$ 的最小奇异值. 则由特征值的分离理论, 有

$$\sigma_n \leq \sigma_n(\mathbf{A} - \alpha \mathbf{I}) + \|\mathbf{E}\|_2 \leq \varepsilon_1 + \varepsilon_2. \quad (7.77)$$

将 $\mathbf{x}^{(0)}$ 按 $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n$ 展开, 有

$$\mathbf{x}^{(0)} = \sum_{i=1}^n \beta_i \mathbf{u}_i, \quad \sum_{i=1}^n |\beta_i|^2 = \|\mathbf{x}^{(0)}\|_2^2 = 1.$$

从而有

$$\begin{aligned} \tilde{\mathbf{y}}^{(1)} &= (\mathbf{A} - \alpha \mathbf{I} + \mathbf{E})^{-1} \mathbf{x}^{(0)} \\ &= \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^H \left(\sum_{i=1}^n \beta_i \mathbf{u}_i \right) = \sum_{i=1}^n \frac{\beta_i}{\sigma_i} \mathbf{v}_i. \end{aligned} \quad (7.78)$$

于是有

$$\|\tilde{\mathbf{y}}^{(1)}\|_2 = \left(\sum_{i=1}^n \left| \frac{\beta_i}{\sigma_i} \right|^2 \right)^{\frac{1}{2}} \geq \frac{|\beta_n|}{\sigma_n} \geq \frac{|\beta_n|}{\varepsilon_1 + \varepsilon_2}.$$

这样一来, 只要 $|\beta_n|$ 不是很小 (即 $\mathbf{x}^{(0)}$ 在 \mathbf{u}_n 方向上不是十分亏损), $\|\tilde{\mathbf{y}}^{(1)}\|_2$ 就会很大. 因此, 通常反幂法只需要迭代一次就足够了.

上述分析表明, 利用反幂法求特征向量时, 位移量 α 取为较精确的近似特征值最好. 此时, 一般只需要迭代一次就可以得到很好的近似特征向量. 因此, 通常总是在用某种方法求得 \mathbf{A} 的近似特征值之后, 再利用反幂法求对应的特征向量. 连同 QR 方法一起来使用反幂法的基本步骤如下.

算法 7.9

步 1, 输入矩阵 $A = (a_{ij}) \in \mathbb{R}^{n \times n}$.

步 2, 上 Hessenberg 化. 用算法 7.4 计算 A 的上 Hessenberg 分解 $A := Q^T A Q$.

步 3, 用双重步位移隐式 QR 方法 (算法 7.7) 求出 A 的特征值, 而不累积正交变换矩阵.

步 4, 对每个计算得到的近似特征值 $\tilde{\lambda}$, 在式 (7.66) 中取位移参数 $\alpha = \tilde{\lambda}$ 进行迭代, 求出特征向量 z .

步 5, 计算 $x = Qz$ (则 x 就是对应于 $\tilde{\lambda}$ 的近似特征向量).

算法 7.9 的 MATLAB 程序如下:

```
function [Lam,V,iter,ki]=ddiqr_eigvec(A,tol)
%用双重步位移隐式QR方法求实方阵的全部特征值和相应的特征向量
%输入:A为n阶实方阵,tol为控制精度(默认是1.e-5)
%输出:Lam为A的全部特征值,V为A的全部特征向量,iter为迭代次数
if nargin<2, tol=1e-5; end
n=size(A,1); x=rand(n,1); %x=ones(n,1);
Lam=zeros(n,1); V=zeros(n);
[A,Q]=mhessen(A); %调用上Hessenberg化程序
%调用双重步位移隐式QR方法求全部特征值
[iter,lambd]=ddiqr_eig(A,tol);
for i=1:n
    [lam,v,k]=mvpower(A,x,lambd(i)); %调用反幂法程序
    V(:,i)=v; ki(i)=k; Lam(i)=lam;
end
V=Q*V; %V的每一列为特征向量
```

例 7.9 用算法 7.9 求例 7.7 中矩阵 A 的全部特征值和特征向量.

解 编写脚本 M 文件 ex79.m, 然后在 MATLAB 命令窗口执行之, 得计算结果:

```
>> ex79
双位移QR方法结果      eig函数计算结果
18.4123 + 0.0000i    18.4123 + 0.0000i
11.1805 + 0.0000i    11.1805 + 0.0000i
 1.7099 - 4.2522i     4.4983 + 0.0000i
 1.7099 + 4.2522i     1.7099 + 4.2522i
-2.2327 + 0.0000i     1.7099 - 4.2522i
 4.4983 + 0.0000i    -2.2327 + 0.0000i
-0.2783 + 0.0000i    -0.2783 + 0.0000i
iter =
```

```

      8
ki =
      2      2      2      2      2      2      2
err =
      1.7959e-05

```

7.5 Givens-Householder 方法

7.3 节中的 Jacobi 方法通过 Givens 变换构造一个正交相似矩阵序列 $\{\mathbf{A}_{k+1}\}$ ($\mathbf{A}_{k+1} = \mathbf{Q}_k^T \mathbf{A}_k \mathbf{Q}_k$), 使得 \mathbf{A}_{k+1} 趋于一个对角矩阵, 从而得到 \mathbf{A} 的全部特征值 λ_i ($i = 1, 2, \dots, n$). Givens 指出, 如果不要求 \mathbf{A}_k 趋于对角矩阵而是三对角矩阵, 则这个过程可以是有限步的. 而 Householder 建议如果 \mathbf{Q}_k 取为形如 $\mathbf{H}_k = \mathbf{I} - 2\mathbf{u}\mathbf{u}^T$ 的 Householder 矩阵, 则可以更有效地实现这个三对角化的过程.

将 \mathbf{A} 化为三对角矩阵 $\mathbf{T} = \mathbf{H}^T \mathbf{A} \mathbf{H}$ 后, 还要求 \mathbf{T} 的特征值. 为此, Givens 根据 $\mathbf{T} - \lambda \mathbf{I}$ 的顺序主子式构成的 Sturm 序列这一事实提出了计算 \mathbf{T} 的特征值的二分法.

为计算 \mathbf{A} 的特征向量, 要先计算 \mathbf{T} 的特征向量. 一个有效的方法是反迭代法.

这种先通过 Householder 变换化对称矩阵 \mathbf{A} 为对称三对角矩阵 \mathbf{T} , 然后通过求 \mathbf{T} 的特征值和特征向量来求得 \mathbf{A} 的特征值和特征向量的方法称为 Givens-Householder 方法. 对称矩阵的三对角化即为上 Hessenberg 化, 这在 7.4.1 节中已经讨论过, 可以执行算法 7.4 得到. 下面只需考虑如何计算对称三对角矩阵 \mathbf{T} 的特征值.

7.5.1 求对称三对角矩阵特征值的二分法

考虑实对称三对角矩阵

$$\mathbf{T} = \begin{bmatrix} \alpha_1 & \beta_2 & & & \\ \beta_2 & \alpha_2 & \beta_3 & & \\ & \ddots & \ddots & \ddots & \\ & & \beta_{n-1} & \alpha_{n-1} & \beta_n \\ & & & \beta_n & \alpha_n \end{bmatrix}.$$

假定 \mathbf{T} 是不可约的, 即 $\beta_i \neq 0$ ($i = 2, \dots, n$). 否则, 可以将 \mathbf{T} 分解成几个阶数更小的不可约三对角矩阵.

令 $p_k(\lambda)$ 表示矩阵 $\mathbf{T} - \lambda \mathbf{I}$ 的 k 阶顺序主子式, 即

$$p_k(\lambda) = \det(\mathbf{T} - \lambda \mathbf{I})_k$$

$$= \det \begin{bmatrix} \alpha_1 - \lambda & \beta_2 & & & \\ \beta_2 & \alpha_2 - \lambda & \beta_3 & & \\ & \ddots & \ddots & \ddots & \\ & & \beta_{k-1} & \alpha_{k-1} - \lambda & \beta_k \\ & & & \beta_k & \alpha_k - \lambda \end{bmatrix}, \quad k = 1, 2, \dots, n.$$

若令 $p_0(\lambda) \equiv 1$, 则有递推式

$$\begin{aligned} p_0(\lambda) &= 1, \quad p_1(\lambda) = \alpha_1 - \lambda, \\ p_{k+1}(\lambda) &= (\alpha_{k+1} - \lambda)p_k(\lambda) - \beta_{k+1}^2 p_{k-1}(\lambda), \quad k = 1, \dots, n-1. \end{aligned} \quad (7.79)$$

容易证明, 多项式序列 $\{p_k(\lambda)\}_0^n$ 具有如下四个特征.

引理 7.2 序列 $\{p_k(\lambda)\}_0^n$ 满足

$$p_k(-\infty) > 0, \quad p_k(+\infty) \begin{cases} > 0, & k \text{ 为偶数,} \\ < 0, & k \text{ 为奇数.} \end{cases} \quad (7.80)$$

证明 首先用归纳法证明

$$p_k(\lambda) = (-1)^k \lambda^k + g_{k-1}(\lambda), \quad (7.81)$$

式中: $g_{k-1}(\lambda)$ 为 λ 的 $k-1$ 次多项式.

事实上, 当 $k=0$ 时, 显然式 (7.81) 成立. 假设 $k-1$ 时结论成立. 则对 k , 有

$$\begin{aligned} p_k(\lambda) &= (\alpha_k - \lambda)p_{k-1}(\lambda) - \beta_k^2 p_{k-2}(\lambda) \\ &= (\alpha_k - \lambda)[(-1)^{k-1} \lambda^{k-1} + g_{k-2}(\lambda)] - \beta_k^2 p_{k-2}(\lambda) \\ &= (-1)^k \lambda^k + (-1)^{k-1} \alpha_k \lambda^{k-1} + (\alpha_k - \lambda)g_{k-2}(\lambda) - \beta_k^2 p_{k-2}(\lambda) \\ &= (-1)^k \lambda^k + g_{k-1}(\lambda). \end{aligned}$$

注意到 k 次多项式 $p_k(\lambda)$ 的符号在 $|\lambda|$ 充分大时决定于其首项的符号, 故由上式立即得到引理的结论. 证毕. \square

引理 7.3 $p_0(\lambda)$ 无零点, 相邻两个多项式 $p_k(\lambda)$ 与 $p_{k+1}(\lambda)$ 无公共零点.

证明 $p_0(\lambda)$ 无零点是显然的. 设 $p_k(\lambda)$ 与 $p_{k+1}(\lambda)$ 有公共零点 $\bar{\lambda}$, 则由递推式 (7.79) 可知, $\bar{\lambda}$ 是 p_{k-1} 的零点 (因 $\beta_{k+1} \neq 0$). 依次类推, 可知 $\bar{\lambda}$ 是 $p_{k-2}(\lambda), \dots, p_1(\lambda), p_0(\lambda)$ 的零点, 与 $p_0(\lambda) = 1$ 矛盾. 证毕. \square

引理 7.4 设 $\bar{\lambda}$ 是 $p_k(\lambda)$ ($0 < k < n$) 的零点, 则 $p_{k-1}(\bar{\lambda})p_{k+1}(\bar{\lambda}) < 0$.

证明 因 $\bar{\lambda}$ 是 p_k 的零点, 由引理 7.3 可知, $p_{k-1}(\bar{\lambda}) \neq 0, p_{k+1}(\bar{\lambda}) \neq 0$. 由式 (7.79) 可知

$$p_{k+1}(\bar{\lambda}) = -\beta_{k+1}^2 p_{k-1}(\bar{\lambda}),$$

故 $p_{k-1}(\bar{\lambda})p_{k+1}(\bar{\lambda}) = -\beta_{k+1}^2 p_{k-1}(\bar{\lambda})^2 < 0$. 证毕. \square

引理 7.5 $p_k(\lambda)$ ($0 < k < n$) 的零点全是单重的, 并且 $p_k(\lambda)$ 的零点把 $p_{k+1}(\lambda)$ 的零点严格地隔离开来.

证明 对 k 使用归纳法. 事实上, 当 $k=1$ 时, $p_1(\lambda) = \alpha_1 - \lambda$, 其零点为 α_1 . 另外, $p_2(\alpha_1) = -\beta_2^2 < 0$, 且由式 (7.80) 可知, $p_2(-\infty) > 0$, $p_2(+\infty) > 0$. 故 $p_2(\lambda)$ 的根必定在区间 $(-\infty, \alpha_1)$ 与 $(\alpha_1, +\infty)$ 内, 即当 $k=1$ 时, 引理结论成立.

假设当 $k=i-1$ 时结论成立, 即 $p_{i-1}(\lambda)$ 和 $p_i(\lambda)$ 的根全是单根且 $p_{i-1}(\lambda)$ 的根把 $p_i(\lambda)$ 的根严格隔开. 记 $\lambda_s^{(l)}$ 为 $p_l(\lambda)$ ($l=i-1, i$) 的第 s 个根. 则有

$$\lambda_1^{(i)} < \lambda_1^{(i-1)} < \lambda_2^{(i)} < \lambda_2^{(i-1)} < \cdots < \lambda_{i-1}^{(i)} < \lambda_{i-1}^{(i-1)} < \lambda_i^{(i)}.$$

由于

$$p_{i-1}(-\infty) > 0, \quad p_{i-1}(\lambda_1^{(i-1)}) = 0, \quad p_{i-1}(\lambda_2^{(i-1)}) = 0, \quad \cdots, \quad p_{i-1}(\lambda_{i-1}^{(i-1)}) = 0,$$

所以

$$p_{i-1}(\lambda_1^{(i)}) > 0, \quad p_{i-1}(\lambda_2^{(i)}) < 0, \quad p_{i-1}(\lambda_3^{(i)}) > 0, \quad \cdots,$$

即 $p_{i-1}(\lambda_s^{(i)})$ 的符号为 $(-1)^{s+1}$. 又因

$$p_{i+1}(\lambda_s^{(i)}) = -\beta_{i+1}^2 p_{i-1}(\lambda_s^{(i)}),$$

故

$$p_{i+1}(-\infty) > 0, \quad p_{i+1}(\lambda_1^{(i)}) < 0, \quad p_{i+1}(\lambda_2^{(i)}) > 0, \quad p_{i+1}(\lambda_3^{(i)}) < 0, \quad \cdots.$$

于是由零点存在定理, $p_{i+1}(\lambda)$ 在区间

$$(-\infty, \lambda_1^{(i)}), (\lambda_1^{(i)}, \lambda_2^{(i)}), \cdots, (\lambda_{i-1}^{(i)}, \lambda_i^{(i)}), (\lambda_i^{(i)}, +\infty)$$

内都有根. 这样的区间共有 $i+1$ 个, 而 $p_{i+1}(\lambda)$ 恰有 $i+1$ 个根, 故在每个区间内有且仅有一个根, 即当 $k=i$ 时结论成立. 由归纳法原理, 引理得证. \square

定义 7.5 对于一个多项式序列 $\{p_k(\lambda)\}_0^n$, 如果它具有引理 7.3 至引理 7.5 的性质, 则称这个多项式序列为 Sturm 序列.

对于 Sturm 序列, 任意取定 α , 则

$$p_0(\alpha), p_1(\alpha), \cdots, p_k(\alpha), \quad k \leq n \quad (7.82)$$

是一个有限序列. 用 $s_k(\alpha)$ 表示式 (7.82) 中每相邻两数符号一致的数目, 并称其为该序列的同号数. 若序列中某一项 $p_k(\alpha) = 0$, 则约定 $p_k(\alpha)$ 的符号与 $p_{k-1}(\alpha)$ 的符号相同. 根据引理 7.3, 此时必有 $p_{k-1}(\alpha) \neq 0$. 例如, 数列 $\{2, 4, 8, 16, -10, -12, 14\}$ 的同号数 $s_6 = 4$, 而数列 $\{-1, -3, -5, 0, 6, 7\}$ 的同号数 $s_5 = 4$.

下面的定理是求对称三对角矩阵特征值二分法的理论基础.

定理 7.16 设 T 是不可约对称三对角矩阵, α 是任意实数, 则 $s_k(\alpha)$ 等于 $p_k(\lambda) = 0$ 在区间 $[\alpha, +\infty)$ 内根的个数.

证明 用归纳法. 当 $k=0$ 时, 由于 $p_0(\alpha)=1$, 故结论显然成立. 当 $k=1$ 时, $p_1(\lambda)=\alpha_1-\lambda$. 若 $p_1(\alpha)<0$, 则 $s_1(\alpha)=0$. 注意到 $p_1(-\infty)>0$, 可知 $p_1(\lambda)=0$ 的根必在 $(-\infty, \alpha)$ 内, 因此在 $[\alpha, +\infty)$ 无根, 即根的个数为 0 个. 若 $p_1(\alpha)\geq 0$, 则 $s_1(\alpha)=1$. 注意到 $p_1(+\infty)<0$, 可知此时 $p_1(\lambda)=0$ 的唯一根必在 $[\alpha, +\infty)$ 内. 故当 $k=1$ 时, 结论为真.

现假定 $k\leq l$ 时结论成立. 记 $s_l(\alpha)=r$, $p_l(\lambda)$ 的零点为

$$\xi_l < \xi_{l-1} < \cdots < \xi_1.$$

由归纳法假设, 有

$$\xi_l < \xi_{l-1} < \cdots < \xi_{r+1} < \alpha \leq \xi_r < \cdots < \xi_1. \quad (7.83)$$

设 $p_{l+1}(\lambda)$ 的零点为

$$\mu_{l+1} < \mu_l < \cdots < \mu_{r+1} < \mu_r < \cdots < \mu_1,$$

则显然有

$$p_l(\alpha) = \prod_{i=1}^l (\xi_i - \alpha), \quad p_{l+1}(\alpha) = \prod_{i=1}^{l+1} (\mu_i - \alpha). \quad (7.84)$$

且据引理 7.5 有 $p_l(\lambda)$ 与 $p_{l+1}(\lambda)$ 的零点相互隔离:

$$\mu_{l+1} < \xi_l < \mu_l < \cdots < \xi_{r+1} < \mu_{r+1} < \xi_r < \mu_r < \cdots < \xi_1 < \mu_1, \quad (7.85)$$

由式 (7.83) 和式 (7.85) 可知, $\xi_{r+1}, \mu_{r+1}, \alpha, \xi_r$ 之间的相互位置只能出现下面四种情况:

(1) $\xi_{r+1} < \alpha < \mu_{r+1}$. 此时 $p_{l+1}(\lambda)$ 在 $[\alpha, +\infty)$ 上恰有 $r+1$ 个零点. 而由式 (7.84) 可知 $p_l(\alpha)$ 的符号为 $(-1)^{l-r}$ 与 $p_{l+1}(\alpha)$ 的符号 $(-1)^{l+1-(r+1)}$ 相同, 故 $s_{l+1}(\alpha) = s_l(\alpha) + 1 = r + 1$.

(2) $\mu_{r+1} < \alpha < \xi_r$. 此时 $p_{l+1}(\lambda)$ 在 $[\alpha, +\infty)$ 上恰有 r 个零点. 而由式 (7.84) 可知 $p_l(\alpha)$ 的符号 $(-1)^{l-r}$ 与 $p_{l+1}(\alpha)$ 的符号 $(-1)^{(l+1)-r}$ 相反, 于是 $s_{l+1}(\alpha) = s_l(\alpha) = r$.

(3) $\mu_{r+1} = \alpha < \xi_r$. 此时 $p_{l+1}(\lambda)$ 在 $[\alpha, +\infty)$ 上恰有 $r+1$ 个零点. 因为此时 $p_{l+1}(\alpha) = 0$, 按约定 $p_{l+1}(\alpha)$ 与 $p_l(\alpha)$ 符号相同, 故仍有 $s_{l+1}(\alpha) = s_l(\alpha) + 1 = r + 1$.

(4) $\mu_{r+1} < \alpha = \xi_r$. 此时 $p_{l+1}(\lambda)$ 在 $[\alpha, +\infty)$ 上恰有 r 个零点. 注意到 $p_l(\alpha) = 0$, 按约定, 它与 $p_{l-1}(\alpha)$ 符号相同. 而由引理 7.4 知, $p_{l+1}(\alpha)$ 与 $p_{l-1}(\alpha)$ 反号, 故 $p_{l+1}(\alpha)$ 与 $p_l(\alpha)$ 反号. 从而 $s_{l+1}(\alpha) = s_l(\alpha) = r$.

这就对 $k=l+1$ 证明了结论成立. 由归纳法原理知定理对一切 $k\leq n$ 都成立. \square

推论 7.4 设 T 是不可约对称三对角矩阵, $p_k(\lambda)$ 是 $T - \lambda I$ 的前 k 阶主子式, $s_n(\alpha)$ 是序列 $\{p_k(\alpha)\}_0^n$ 的同号数, 则 $s_n(\alpha)$ 等于 T 在区间 $[\alpha, +\infty)$ 中特征值的个数.

利用推论 7.4, 可以确定在任意区间 $(\alpha, \beta]$ 中所含 T 的特征值的数目. 实际上, 它就等于 $s_n(\alpha) - s_n(\beta)$. 而且若 T 的特征值排列为

$$\lambda_n < \lambda_{n-1} < \cdots < \lambda_m < \cdots < \lambda_1,$$

$s_n(\alpha) \geq m > s_n(\beta)$, 则 $\lambda_m \in [\alpha, \beta)$.

综上所述, 可以给出计算 T 的特征值 λ_m 的二分法.

算法 7.10 (计算三对角矩阵特征值的二分法)

- 步 1, 输入包含 λ_m 的初始区间 $[a_0, b_0]$ (如取 $a_0 = -\|T\|, b_0 = \|T\|$). 置 $k := 0$.
- 步 2, 计算 $[a_k, b_k]$ 的中点 $c_k = \frac{a_k + b_k}{2}$ 及 $\{p_i(c_k)\}_{i=0}^n$ 的同号数 $s_n(c_k)$.
- 步 3, 若 $s_n(c_k) \geq m$, 则 $a_{k+1} := c_k, b_{k+1} := b_k$; 否则, $a_{k+1} := a_k, b_{k+1} := c_k$.
- 步 4, 若 $|b_{k+1} - a_{k+1}| \leq \varepsilon$, 则令 $\lambda_m \approx \frac{1}{2}(a_{k+1} + b_{k+1})$, 停算. 否则, 置 $k := k + 1$, 转步 2.

在算法 7.10 中, λ_m 始终属于区间 $[a_k, b_k]$ ($k = 0, 1, \dots$). 当 k 充分大时, $[a_k, b_k]$ 的长度 $\frac{b_0 - a_0}{2^k}$ 可以小于任意指定的精度, 此时区间中点 c_k 作为 λ_m 的近似值, 其误差不会超过 $\varepsilon/2$. 若二分法的计算是精确的, 则可以得到任意指定精度的近似特征值. 但实际计算中会有舍入误差的影响. 虽然如此, 利用后验误差分析的方法可以证明二分法是一个数值稳定的方法.

注 7.4 在算法 7.10 中, 每一步都要计算 n 个多项式 $p_i(\lambda)$ ($i = 1, 2, \dots, n$) 在区间 $[a_k, b_k]$ 的中点 c_k 的值 (因 $p_0(\lambda) = 1$, 故不需要计算), 这很容易发生“上溢”或“下溢”现象.

下面讨论 $s = s_n(\alpha)$ 的计算. 由于

$$p_k(\alpha) = (\alpha_k - \alpha)p_{k-1}(\alpha) - \beta_k^2 p_{k-2}(\alpha), \quad k = 1, 2, \dots, n,$$

这里 $p_0(\alpha) = 1, \beta_1 = 0, p_{-1}(\alpha) = 0$. 令 $s = 0$. 若 $p_{k-1}(\alpha)p_k(\alpha) > 0$, s 加 1. 若 $p_{k-1}(\alpha)p_k(\alpha) < 0$, s 不变. 若 $p_{k-1}(\alpha)p_k(\alpha) = 0$, 由于 $p_{k-1}(\alpha)$ 与 $p_k(\alpha)$ 不能同时为 0, 分为两种情况. 第一种情况是 $p_{k-1}(\alpha) = 0, p_k(\alpha) \neq 0$, 此时 $p_{k-1}(\alpha)$ 与 $p_{k-2}(\alpha) (\neq 0)$ 同号. 若 $p_{k-2}(\alpha)p_k(\alpha) > 0$, 则 s 加 1, 否则 s 不变. 第二种情况是 $p_{k-1}(\alpha) \neq 0, p_k(\alpha) = 0$. 此时 $p_k(\alpha)$ 与 $p_{k-1}(\alpha) (\neq 0)$ 同号, s 加 1.

因此, 如果引入

$$q_k(\alpha) = p_k(\alpha)/p_{k-1}(\alpha), \quad k = 1, 2, \dots, n,$$

则

$$q_k(\alpha) = (\alpha_k - \alpha) - \beta_k^2/q_{k-1}(\alpha), \quad k = 1, 2, \dots, n. \quad (7.86)$$

对于任意的 $1 \leq k \leq n$, ① 若 $q_{k-1}(\alpha) \neq 0$, 当式 (7.86) 的计算结果 $q_k(\alpha) \geq 0$ 时, s 加 1, 否则, s 不变; ② 若 $q_{k-1}(\alpha) = 0$, 这表明 $p_{k-1}(\alpha) = 0$, 从而有 $p_k(\alpha) = -\beta_k^2 p_{k-2}(\alpha)$. 由于 $\beta_k \neq 0, p_k(\alpha)$ 与 $p_{k-2}(\alpha)$ 异号, 故 $p_k(\alpha)$ 与 $p_{k-1}(\alpha)$ 同号, s 加 1.

为此, 定义一个新的序列 $\{q_k(\lambda)\}_{k=1}^n$ 如下:

$$q_1(\lambda) = \alpha_1 - \lambda,$$

$$q_k(\lambda) = \begin{cases} \alpha_k - \lambda - \frac{\beta_k^2}{q_{k-1}(\lambda)}, & q_{k-1}(\lambda) \neq 0, \\ 1, & q_{k-1}(\lambda) = 0, \end{cases}$$

$$k = 2, 3, \dots, n.$$

则 $\{p_k(\lambda)\}_{k=0}^n$ 的同号数 $s_n(\lambda)$ 等于 $\{q_k(\lambda)\}_{k=1}^n$ 中非负数的数目. 故在二分法的实际计算中, 总是通过计算 $\{q_k(\lambda)\}_{k=1}^n$ 的非负数的数目来代替计算 $s_n(\lambda)$.

7.5.2 二分法的程序实现

本节考虑算法 7.10 的程序实现. 首先考虑多项式序列 $\{p_k(\lambda)\}_{k=1}^n$ 和 $\{q_k(\lambda)\}_{k=1}^n$. 编制 MATLAB 程序如下:

```
%本函数计算Sturm序列{p_k(t)}的值
function p=pkfun(t,a,b)
n=length(a);
p(1)=a(1)-t; p(2)=(a(2)-t)*p(1)-b(1);
for k=3:n
    p(k)=(a(k)-t)*p(k-1)-b(k-1)^2*p(k-2);
end
```

%本函数计算防止溢出的多项式序列{q_k(t)}的值

```
function q=qkfun(t,a,b)
n=length(a);
q(1)=a(1)-t;
for k=2:n
    if q(k-1)==0
        q(k)=1;
    else
        q(k)=(a(k)-t-b(k-1)^2/q(k-1));
    end
end
```

用数据

$$a = [0.8147, 0.9058, 0.1270, 0.9134, 0.6324, 0.0975, 0.2785, 0.5469]^T,$$

$$b = [0.9575, 0.9649, 0.1576, 0.9706, 0.9572, 0.4854, 0.8003]^T, \quad t = 6.8,$$

测试两个程序, 得

```
>> a=[0.8147,0.9058,0.1270,0.9134,0.6324,0.0975,0.2785,0.5469]';
>> b=[0.9575,0.9649,0.1576,0.9706,0.9572,0.4854,0.8003]';
>> p=pkfun(6.8,a,b)
p =
    1.0e+06 *
    -0.0000    0.0000   -0.0002    0.0013   -0.0079    0.0517   -0.3355    2.0646
>> q=qkfun(6.8,a,b)
q =
   -5.9853   -5.7410   -6.5108   -5.8828   -6.0075   -6.5500   -6.4855   -6.1543
```

可以看出序列 $\{q_k(\lambda)\}_{k=1}^n$ 的计算的确能有效地防止“溢出”。

现在结合算法 7.10 和注 7.4, 编制 MATLAB 程序如下:

```
%Givens-Householder方法程序-givens_househ.m
function [lambda,k]=givens_househ(T,m,tol,max_it)
%用Givens-Householder求实对称三对角矩阵的第m个特征值
%输入:T为n阶实对称方阵,m特征值序号,
%    tol为容许误差,max_it为最大迭代次数
%输出:lambda为第m个特征值,k为满足精度迭代次数
if nargin<4, max_it=500; end
if nargin<3, tol=1e-6; end
n=size(T,1); k=0;
b=norm(T,inf); a=-b;
alpha=diag(T); beta=diag(T,1);
while(k<max_it)
    k=k+1; c=(a+b)/2;
    q=qkfun(c,alpha,beta);
    tt=find(q>=0);
    sn=length(tt);
    if(sn>=m), a=c; else, b=c; end
    if abs(b-a)<=tol,
        lambda=c; break;
    end
end
end
```

下面看一个计算实例。

例 7.10 用 Givens-Householder 方法计算矩阵

$$A = \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & 2 & 2 & \cdots & 2 \\ 1 & 2 & 3 & \cdots & 3 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 2 & 3 & \cdots & n \end{bmatrix}$$

的全部特征值。取 $n = 12$ 。

解 编制 MATLAB 脚本文件 ex710.m, 然后在命令窗口执行之, 即可得到 A 全部特征值以及计算时间。

7.5.3 特征向量的计算

关于 A 的特征向量的计算, 既可以直接从 A 出发借助已求得的近似特征值利用反迭代法计算, 也可以先用反迭代法求出 T 的相应特征向量, 然后再利用三对角化 (上

Hessenberg 化) 过程所得到的变换矩阵 H_1, H_2, \dots, H_{n-2} 将其变换为 A 的特征向量. 后者由于其方法简单、计算量小而得到推荐.

设 (λ, x) 是矩阵 T 的一个特征对, 即 $Tx = \lambda x$. 由于

$$T = (H_1 H_2 \cdots H_{n-2})^T A (H_1 H_2 \cdots H_{n-2}),$$

故

$$A(H_1 H_2 \cdots H_{n-2})x = \lambda(H_1 H_2 \cdots H_{n-2})x,$$

于是

$$z = (H_1 H_2 \cdots H_{n-2})x \quad (7.87)$$

即为 A 的相应于 λ 的特征向量.

实际计算时不是直接按式 (7.87) 进行 $n-2$ 个矩阵乘向量的运算, 而是利用 Householder 矩阵 $H_k = I - 2u_k u_k^T$ 的特殊形式按下列公式计算:

$$\begin{cases} y_{n-1} = x, \\ y_k = H_k y_{k+1} = y_{k+1} - 2(u_k^T y_{k+1})u_k, \\ z = y_1. \end{cases}$$

因此, 问题归结为求不可约对称三对角矩阵 T 的特征向量 x . 通常采用所谓的反迭代法. 反迭代法本质上是用原点位移的反幂法求矩阵的特征向量, 其基本迭代格式如下:

$$\begin{cases} (T - \mu I)v_k = z_{k-1}, \\ z_k = v_k / \|v_k\|_2, \quad k = 1, 2, \dots, \end{cases} \quad (7.88)$$

式中: μ 为事先取定的位移参数; z_0 为初始向量.

通常用 Givens-Householder 方法 (或其他方法, 如 QR 方法等) 求得某个特征值的近似值后, 即用其作为位移参数, 然后再用式 (7.88) 进行迭代, 会有很快的收敛速度. 反迭代法的基本步骤如下.

算法 7.11 (计算特征向量的反迭代法)

步 1, 计算 A 的上 Hessenberg 分解: $Q^T A Q = T$.

步 2, 使用算法 7.10 计算 T 的近似特征值 $\bar{\lambda}_1 > \bar{\lambda}_2 > \dots > \bar{\lambda}_n$.

步 3, 对每个计算得到的特征值 $\bar{\lambda}_i (i = 1, 2, \dots, n)$, 在式 (7.88) 中取 $\mu := \bar{\lambda}_i$ 进行迭代, 求出特征向量 z_i .

步 4, 计算 $x_i = Q z_i$ (则 x_i 就是对应于 $\bar{\lambda}_i$ 的近似特征向量).

Givens-Householder 方法可以用来计算对称矩阵的全部特征值问题. 但它可以灵活地求解各种部分特征值问题, 如求若干最大 (最小) 特征值, 求位于给定区间 $[\alpha, \beta]$ 上的特征值等. 这一点它比其他部分特征值问题的算法显得更为优越, 所以这一方法通常也用来计算对称矩阵的部分特征值问题.

算法 7.11 的 MATLAB 程序如下:

```

%基于Givens-Householder方法的求特征向量的反迭代法程序-ghvector.m
function [Lam,V,ki]=ghvector(A,tol,max_it)
%用反迭代法求实对称矩阵A的第m大的特征向量
%输入:A为n阶对称方阵,m为按降幂排列特征值序号,
%      tol为容许误差,max_it为最大迭代次数
%输出:lambda为第m大的特征值,x为对应的特征向量,k为迭代次数
n=size(A,1); Lam=zeros(n,1); x=rand(n,1); %x=ones(n,1);
[T,Q]=mhessen(A); %调用上Hessenberg化程序
for i=1:n
    [la]=givens_househ(T,i,tol,max_it); Lam(i)=la;
    [lam,v,k2]=mvpower(T,x,Lam(i)); %调用反幂法程序
    Lam(i)=lam; V(:,i)=v; ki(i)=k2;
    norm(T*v-Lam(i)*v),
end
V=Q*V; %V的每一列为特征向量

```

例 7.11 用算法 7.11 计算例 7.11 中矩阵 A 的全部特征值和特征向量. 取 $n = 15$.

解 编制 MATLAB 脚本文件 ex711.m, 然后在命令窗口执行之, 即可得到 A 全部特征值和特征向量以及计算时间.

7.6 Krylov 子空间方法

由于目前计算机存储空间和运算速度的限制, QR 方法以及由其导出的各类方法只适用于小型矩阵特征值和特征向量的计算, 而并不能用于求解大型矩阵的特征值问题. 对于大型稀疏矩阵的特征值问题而言, 目前较为有效和实用的一种求解方法就是 Krylov 子空间方法. 特别是近几年, 随着重新开始技术的不断完善, 使得这类方法的效率越来越高, 适用范围越来越广. 本节介绍与 Krylov 子空间方法有关的一些基本概念和重要理论, 并且详细地阐述目前较为成熟的 Arnoldi 方法和 Lanczos 方法.

设 $A \in \mathbb{R}^{n \times n}$ 是一个大型稀疏的矩阵, 即 $n \gg 1$ (如 $n \geq 10^6$), 而且 A 的非零元很少 (如不超过 10%). 对于这样矩阵的特征值问题, 一般来讲用 QR 方法来求解已经是不可能的事情. 例如, 假如希望用 QR 方法来计算 A 的特征值, 则完成这一计算任务所需的运算量约为 $10n^3$. 如果 $n = 10^7$, 而且假定用一个万亿次的计算机来执行这一计算任务, 则所需的计算时间为

$$\frac{10 \times 10^{21}}{10^{12} \times 3600 \times 24 \times 365} \approx 317 \text{ (年)}$$

此外, QR 方法所需的存储量是 $O(n^2) = O(10^{14}) = O(15G)$, 这也是现在的计算机上无法实现的事情.

那么如何来求解大型的矩阵计算问题呢? 一个直观而朴素的想法是: 先用一个小型的问题“逼近”所考虑的大型问题, 然后用这个小型问题的解近似所求的解. 实现这

一想法的一个途径是, 先选择一个适当的 k 维子空间 $\mathcal{X}_k \subset \mathbb{R}^n, k \ll n$, 然后将所考虑的问题在某种特定的意义下“投影”到 \mathcal{X}_k 上变成一个小型的问题.

目前, 常用的子空间 \mathcal{X}_k 是 Krylov 子空间:

$$\mathcal{K}_k(\mathbf{A}, \mathbf{v}) = \text{span}\{\mathbf{v}, \mathbf{A}\mathbf{v}, \dots, \mathbf{A}^{k-1}\mathbf{v}\},$$

式中: $\mathbf{v} \in \mathbb{R}^n$ 是人为给定的.

选择这一空间是因为它有一些特有的优点:

- (1) 所需信息量少: 只需一个矩阵 \mathbf{A} , 一个向量 \mathbf{v} 和一个正整数 k .
- (2) 计算简单: 只涉及矩阵乘向量, 可以充分利用 \mathbf{A} 所具有的稀疏性.
- (3) 容易扩张: $\mathcal{K}_{k+1}(\mathbf{A}, \mathbf{v}) = \mathcal{K}_k(\mathbf{A}, \mathbf{v}) + \text{span}\{\mathbf{A}^k \mathbf{v}\}$.

当所有的子空间确定之后, 下一步的任务就是如何将所考虑的大型问题“投影”到 \mathcal{X}_k 上. 设 $\mathbf{V}_k \in \mathbb{R}^{n \times k}$ 是 \mathcal{X}_k 上的一组标准正交基构成的矩阵, 即若记 $\mathbf{V}_k = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k]$, 则有

$$\mathcal{X}_k = \text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}, \quad \mathbf{V}_k^T \mathbf{V}_k = \mathbf{I}. \quad (7.89)$$

有了 \mathbf{V}_k 之后, 就可以计算 \mathbf{A} 在 \mathcal{X}_k 上的投影

$$\mathbf{H}_k = \mathbf{V}_k^T \mathbf{A} \mathbf{V}_k.$$

现在来看 \mathbf{H}_k 的几何解释. 令 $\mathcal{A}: \mathbb{R}^n \mapsto \mathbb{R}^n$ 为

$$\mathbf{x} \mapsto \mathbf{A}\mathbf{x}, \quad \forall \mathbf{x} \in \mathbb{R}^n.$$

则有

$$\begin{array}{ccc} \mathcal{X}_k & \xrightarrow{\mathcal{A}|_{\mathcal{X}_k}} & \mathbb{R}^n \\ & \searrow & \downarrow \mathcal{P}_{\mathcal{X}_k} \\ & \mathcal{P}_{\mathcal{X}_k} \circ \mathcal{A}|_{\mathcal{X}_k} & \mathcal{X}_k \end{array}$$

而且 \mathbf{H}_k 正好是线性算子 $\mathcal{P}_{\mathcal{X}_k} \circ \mathcal{A}|_{\mathcal{X}_k}$ 在基向量 $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ 下的矩阵表示, 其中 $\mathcal{A}|_{\mathcal{X}_k}$ 表示算子 \mathcal{A} 在子空间 \mathcal{X}_k 上的限制, $\mathcal{P}_{\mathcal{X}_k}$ 表示 \mathcal{X}_k 上的正交投影算子, $\mathcal{P}_{\mathcal{X}_k} \circ \mathcal{A}|_{\mathcal{X}_k}$ 表示算子 $\mathcal{P}_{\mathcal{X}_k}$ 和 $\mathcal{A}|_{\mathcal{X}_k}$ 的复合.

事实上, 只要注意到: ① $\mathbf{x} \in \mathcal{X}_k$ 当且仅当存在 $\mathbf{y} \in \mathbb{R}^k$ 使得 $\mathbf{x} = \mathbf{V}_k \mathbf{y}$; ② $\mathcal{P}_{\mathcal{X}_k}$ 的矩阵表示为 $\mathbf{V}_k \mathbf{V}_k^T$. 有

$$\mathcal{P}_{\mathcal{X}_k} \circ \mathcal{A}|_{\mathcal{X}_k}(\mathbf{V}_k \mathbf{y}) = \mathbf{V}_k \mathbf{V}_k^T \mathbf{A} \mathbf{V}_k \mathbf{y}, \quad \forall \mathbf{y} \in \mathbb{R}^k,$$

即

$$\mathcal{P}_{\mathcal{X}_k} \circ \mathcal{A}|_{\mathcal{X}_k} \mathbf{V}_k = \mathbf{V}_k \mathbf{H}_k.$$

这表明, \mathbf{H}_k 正好是线性算子 $\mathcal{A}|_{\mathcal{X}_k}$ 在基向量 $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ 下的矩阵表示.

只要有了 \mathbf{H}_k 和 \mathbf{V}_k , 涉及大矩阵 \mathbf{A} 的计算问题就可用相应的小矩阵 \mathbf{H}_k 对应的问题来解决.

设 $(\lambda, \boldsymbol{x})$ 是矩阵 \boldsymbol{A} 的一个特征对, 即

$$\boldsymbol{A}\boldsymbol{x} = \lambda\boldsymbol{x}. \quad (7.90)$$

再假定

$$\boldsymbol{x} = \boldsymbol{V}_k\boldsymbol{y} + \boldsymbol{v} \approx \boldsymbol{V}_k\boldsymbol{y}.$$

于是, 由式 (7.90), 有

$$\boldsymbol{A}\boldsymbol{V}_k\boldsymbol{y} \approx \lambda\boldsymbol{V}_k\boldsymbol{y}.$$

从而有

$$\boldsymbol{H}_k\boldsymbol{y} = \boldsymbol{V}_k^T\boldsymbol{A}\boldsymbol{V}_k\boldsymbol{y} \approx \lambda\boldsymbol{y}.$$

因此, 若 $(\lambda, \boldsymbol{y})$ 是 $\boldsymbol{H}_k = \boldsymbol{V}_k^T\boldsymbol{A}\boldsymbol{V}_k$ 的特征对, 则可用 $(\lambda, \boldsymbol{V}_k\boldsymbol{y})$ 近似 \boldsymbol{A} 的特征对. 这样就将大型特征值问题 (7.90) 转化为一个小型特征值问题:

$$\boldsymbol{H}_k\boldsymbol{y} = \lambda\boldsymbol{y}.$$

这就是用 Krylov 子空间方法求解大规模特征值问题的基本思想.

7.6.1 Rayleigh-Ritz 投影方法

基于 Krylov 子空间来设计特征值问题数值解法的最基本的技术就是投影. 下面介绍 Rayleigh-Ritz 投影方法及其基本性质.

设 $\boldsymbol{A} \in \mathbb{R}^{n \times n}$, $\mathcal{V} \subset \mathbb{C}^n$ 是一个 m 维子空间 (不一定是 Krylov 子空间), 借助 \mathcal{V} 给出 \boldsymbol{A} 的某些近似特征对. 达到这一目的的一种途径就是正交投影法, 其基本思想是选择 $\mu \in \mathbb{C}$ 和 $\boldsymbol{u} \in \mathcal{V}$ 使得

$$(\boldsymbol{A}\boldsymbol{u} - \mu\boldsymbol{u}) \perp \mathcal{V}. \quad (7.91)$$

通常称条件 (7.91) 为 Galerkin 条件. 这一条件可以等价地表述为

$$(\boldsymbol{A}\boldsymbol{u} - \mu\boldsymbol{u}, \boldsymbol{v}) = 0, \quad \forall \boldsymbol{v} \in \mathcal{V}. \quad (7.92)$$

设 \mathcal{V} 的一组标准正交基为 $\boldsymbol{v}_1, \boldsymbol{v}_2, \dots, \boldsymbol{v}_m$, 并令 $\boldsymbol{V}_m = [\boldsymbol{v}_1, \boldsymbol{v}_2, \dots, \boldsymbol{v}_m]$, 则 \boldsymbol{u} 可表示为 $\boldsymbol{u} = \boldsymbol{V}_m\boldsymbol{y}$, 其中 $\boldsymbol{y} \in \mathbb{C}^m$. 这样式 (7.92) 可以等价地表示为

$$\boldsymbol{V}_m^H(\boldsymbol{A}\boldsymbol{V}_m\boldsymbol{y} - \mu\boldsymbol{V}_m\boldsymbol{y}) = \mathbf{0},$$

即

$$\boldsymbol{A}_m\boldsymbol{y} = \mu\boldsymbol{y}, \quad (7.93)$$

其中 $\boldsymbol{A}_m = \boldsymbol{V}_m^H\boldsymbol{A}\boldsymbol{V}_m$ 正好是 \boldsymbol{A} 关于 \boldsymbol{V}_m 的 Rayleigh 商. 这里, 对 \boldsymbol{A}_m 的每一个特征对 (μ, \boldsymbol{y}) , 称 μ 为 Ritz 值, 而称 $\boldsymbol{u} = \boldsymbol{V}_m\boldsymbol{y}$ 为 Ritz 向量, 称 $(\mu, \boldsymbol{V}_m\boldsymbol{y})$ 为 Ritz 对. 在一定条件下, Ritz 对 $(\mu, \boldsymbol{V}_m\boldsymbol{y})$ 将是 \boldsymbol{A} 的一个很好的近似特征对.

上述的这一求 \boldsymbol{A} 的某些近似特征值和近似特征向量的方法就称为 Rayleigh-Ritz 投影方法, 其基本步骤如下.

算法 7.12 (Rayleigh–Ritz 投影方法)

- 步 1, 求 \mathcal{V} 的一组标准正交基 $\{v_1, v_2, \dots, v_m\}$, 并记 $V_m = [v_1, v_2, \dots, v_m]$.
 步 2, 计算 Rayleigh 商 $A_m = V_m^H A V_m$.
 步 3, 计算 A_m 的特征值, 并选择其中若干所需的作为 A 的近似值: $\mu_1, \mu_2, \dots, \mu_k$.
 步 4, 计算 μ_i 所对应的特征向量 y_i , 并形成 Ritz 向量 $u_i = V_m y_i, i = 1, 2, \dots, k$.

一般来讲, 有 $m \ll n$. 因此, A_m 将是一个阶数很小的矩阵. 这样, 就有很多成熟的数值方法用来完成上述方法中步 3 和步 4 的计算任务. 例如, 可先用 QR 方法来计算 A_m 特征值, 然后再用反幂法来求其对应的特征向量. 另外需要说明的一点是, 由于特征向量比 Schur 向量对扰动敏感得多, 因此, 在某些情况下, 步 4 中的“求 Ritz 向量”可以用“求 Schur 向量”来代替 (所谓 Schur 向量是指实矩阵 $A \in \mathbb{R}^{n \times n}$ 的实 Schur 分解 $A = U^T T U$ 中正交矩阵 U 的列向量, 其中 T 为对角块为 1×1 或 2×2 子矩阵的拟上三角矩阵).

下面给出用 Ritz 对来近似特征对的一个理论依据. 这里假定 $A \in \mathbb{R}^{n \times n}$ 且 $A^T = A$, 即 A 是 n 阶实对称矩阵.

设 $V = [V_k, V_u]$ 是一个 n 阶正交矩阵, 其中 $V_k \in \mathbb{R}^{n \times k}, V_u \in \mathbb{R}^{n \times (n-k)}$. 在实际应用中, 通常 V_k 是已知的, 是由 Lanczos 方法产生的 Krylov 子空间的一组标准正交基构成的, 而 V_u 通常是不知道的. 在下面的讨论中, 将使用如下的符号:

$$T = V^T A V = \begin{bmatrix} V_k^T A V_k & V_k^T A V_u \\ V_u^T A V_k & V_u^T A V_u \end{bmatrix} := \begin{bmatrix} T_k & T_{uk} \\ T_{ku} & T_u \end{bmatrix}. \quad (7.94)$$

Ritz 对作为 A 的特征对的“最佳”逼近有多种理由. 其中之一可以从上面的矩阵看出, 由于 T_{uk} 和 T_u 是不知道的, 而 T_k 是唯一知道的, 因此自然只能借助 T_k 的特征对来构造出 A 的近似特征对. 另一个理由是下面的定理 7.17 所描述的 Rayleigh 商的最佳逼近性.

如所周知, 若对某个 k 阶矩阵 R 有 $AV_k = V_k R$, 则 $\mathcal{R}(V_k)$ 就是 A 的一个不变子空间, 从而 R 的特征值全都是 A 的特征值, 即 $\lambda(R) \subset \lambda(A)$. 如果上述等式不能成立, 自然希望找到一个 k 阶矩阵 R , 使得

$$\|AV_k - V_k R\|_2 = \min,$$

然后用 R 的特征值来近似 A 的特征值. 下面的定理表明 Rayleigh 商 $T_k = V_k^T A V_k$ 就是满足此条件的矩阵.

定理 7.17 设 $A \in \mathbb{R}^{n \times n}$ 满足 $A^T = A$, $V_k \in \mathbb{R}^{n \times k}$ 满足 $V_k^T V_k = I$, 则有

(1) T_k 满足

$$\begin{aligned} \|AV_k - V_k T_k\|_2 &= \min\{\|AV_k - V_k S\|_2 : S^T = S \in \mathbb{R}^{k \times k}\} \\ &= \|T_{ku}\|_2. \end{aligned} \quad (7.95)$$

(2) 若 T_k 的谱分解为 $T_k = YMY^T$, 其中 Y 是正交矩阵, M 是对角矩阵, 则

$$\|A(V_k Y) - (V_k Y)M\|_2 = \min\{\|AP_k - P_k D\|_2 : P_k \in \mathcal{P}, D \in \mathcal{D}\}$$

$$= \|T_{ku}\|_2, \quad (7.96)$$

式中:

$$\mathcal{P} = \{P_k \in \mathbb{R}^{n \times k} : P_k^T P_k = I, \mathcal{R}(P_k) = \mathcal{R}(V_k)\},$$

$$\mathcal{D} = \{D \in \mathbb{R}^{k \times k} : D \text{ 是对角矩阵}\}.$$

证明 (1) 对任意的 $S^T = S \in \mathbb{R}^{k \times k}$, 有

$$\begin{aligned} \|AV_k - V_k S\|_2 &= \|V^T(AV_k - V_k S)\|_2 = \|[V_k, V_u]^T(AV_k - V_k S)\|_2 \\ &= \left\| \begin{bmatrix} V_k^T AV_k - V_k^T V_k S \\ V_u^T AV_k - V_u^T V_k S \end{bmatrix} \right\|_2 = \left\| \begin{bmatrix} T_k - S \\ T_{ku} \end{bmatrix} \right\|_2 \\ &\geq \|T_{ku}\|_2, \end{aligned}$$

而且当 $S = T_k$ 时上述不等式等号成立. 因此, 定理的结论 (1) 成立.

(2) 任意取 $P_k \in \mathcal{P}$ 和 $D \in \mathcal{D}$, 由 \mathcal{P} 的定义可知, 必存在一个 k 阶正交矩阵 U 使得 $P_k = V_k U$. 再由结论 (1) 所证和谱范数的酉不变性, 有

$$\begin{aligned} \|AP_k - P_k D\|_2 &= \|AV_k U - V_k \bar{U} D\|_2 = \|AV_k - V_k U D U^T\|_2 \\ &\geq \|AV_k - V_k T_k\|_2 = \|AV_k - V_k Y M Y^T\|_2 \\ &= \|A(V_k Y) - (V_k Y) M\|_2 = \|T_{ku}\|_2. \end{aligned}$$

注意到 P_k 和 D 的任意性, 可知式 (7.96) 成立. 证毕. \square

注 7.5 从定理 7.17 的证明可以看出, 将谱范数换为 Frobenius 范数定理的结论仍然成立. 定理的结论 (1) 说明, 对于实对称矩阵 A 来说, Rayleigh 商 $T_k = V_k^T A V_k$ 具有最佳逼近性; 而定理的结论 (2) 又说明用 Ritz 对作为其近似特征对是最优的.

一般来讲, T_{ku} 是未知的, 但如果 V_k 是由 Lanczos 过程 (即算法 2.13) 产生的, 则有

$$T = \begin{bmatrix} T_k & T_{uk} \\ T_{ku} & T_u \end{bmatrix} = \left[\begin{array}{ccc|ccc} \alpha_1 & \beta_1 & & & & \\ \beta_1 & \ddots & \ddots & & & \\ & \ddots & \ddots & \beta_{k-1} & & \\ & & & \beta_{k-1} & \alpha_k & \beta_k \\ \hline & & & \beta_k & \alpha_{k+1} & \beta_{k+1} \\ & & & & \beta_{k+1} & \ddots & \ddots \\ & & & & & \ddots & \ddots & \beta_{n-1} \\ & & & & & & \beta_{n-1} & \alpha_n \end{array} \right], \quad (7.97)$$

从而 $T_{ku} = \beta_k e_1 e_k^T$, 其中 β_k 已由算法 2.13 算出, 于是此时有 $\|T_{ku}\|_2 = \beta_k$. 此外, 利用 T_{ku} 还可以给出 Ritz 对作为近似特征对的绝对误差.

定理 7.18 假设符号同前, 并设 T_k 的谱分解为

$$T_k = YMY^T,$$

式中: $Y = [y_1, y_2, \dots, y_k]$ 为正交矩阵; $M = \text{diag}(\mu_1, \mu_2, \dots, \mu_k)$. 则有

(1) 存在 A 的 k 个特征值 $\lambda_1, \dots, \lambda_k$, 使得

$$|\lambda_i - \mu_i| \leq \|T_{ku}\|_2, \quad i = 1, 2, \dots, k. \quad (7.98)$$

(2) 令 $z_i = V_k y_i$, 则

$$\|Az_i - \mu_i z_i\|_2 = \|T_{ku} y_i\|_2, \quad i = 1, 2, \dots, k. \quad (7.99)$$

(3) 若 V_k 是由 Lanczos 过程 (算法 2.13) 得到的, 则有

$$\|T_{ku}\|_2 = |\beta_k|, \quad \|T_{ku} y_i\|_2 = |\beta_k (y_i)_k|, \quad (7.100)$$

式中: $(y_i)_k$ 为 y_i 的最后一个分量.

证明 (1) 令 $\hat{T} = \text{diag}(T_k, T_u)$, 则

$$\|T - \hat{T}\|_2 = \left\| \begin{bmatrix} O & T_{uk} \\ T_{ku} & O \end{bmatrix} \right\|_2 = \|T_{ku}\|_2,$$

然后应用 Weyl 定理即知式 (7.98) 成立.

(2) 直接计算, 有

$$\begin{aligned} \|Az_i - \mu_i z_i\|_2 &= \|V^T(AV_k y_i) - \mu_i V^T(V_k y_i)\|_2 \\ &= \left\| \begin{bmatrix} T_k y_i \\ T_{ku} y_i \end{bmatrix} - \begin{bmatrix} \mu_i y_i \\ 0 \end{bmatrix} \right\|_2 = \left\| \begin{bmatrix} 0 \\ T_{ku} y_i \end{bmatrix} \right\|_2 = \|T_{ku} y_i\|_2, \end{aligned}$$

故式 (7.99) 成立.

(3) 由式 (7.97) 可得式 (7.100). 证毕. \square

7.6.2 Lanczos 方法

Lanczos 方法是日前计算大型稀疏实对称矩阵少数几个特征值和对应特征向量最常用的 Krylov 子空间方法之一. 本节介绍 Lanczos 方法的详细算法及其相关理论.

1. 经典 Lanczos 方法

设 $A^T = A \in \mathbb{R}^{n \times n}$ 已经给定, 希望求 A 的几个最大特征对或最小特征对. 将算法 2.13 与 7.6.1 节所介绍的 Rayleigh-Ritz 方法相结合就可得到完成这一计算任务的经典 Lanczos 方法.

算法 7.13 (经典 Lanczos 方法) 给定一个 n 阶实对称矩阵 A . 本算法计算 A 的少数几个两端特征对.

步 1, 选择初始向量 v , 并令 $v_1 = v/\|v\|_2$.

步 2, 应用算法 2.13 产生一个长度为 k 的 Lanczos 分解:

$$AV_k = V_k T_k + \beta_k v_{k+1} e_k^T.$$

步 3, 计算 T_k 的特征值, 并选择其中满足要求的记为 $\mu_i, i = 1, 2, \dots, \ell$.

步 4, 计算 T_k 对应于 μ_i 的特征向量 y_i , 并形成 Ritz 向量 $u_i = V_k y_i, i = 1, 2, \dots, \ell$.

步 5, 检验是否已经满足要求. 如果不满足要求, 则增加 k , 再返回到步 2, 重复以上各步.

注 7.6 定理 7.18 说明, 若 Rayleigh 商 T_k 的一个特征对 (μ, y) 使得 $|\beta_k| \cdot |e_k^T y|$ 很小, 则 μ 就是 A 的某个特征值的很好近似. 因此, 通常在算法 7.13 中就是以 $|\beta_k| \cdot |e_k^T y|$ 是否已经足够小来判断 Ritz 对 $(\mu, V_k y)$ 是否可以作为 A 的近似特征对.

经典 Lanczos 方法的 MATLAB 程序如下:

```
function [mu,U,lambda]=Class_Lanczos(A,v,k)
%输入:A为对称矩阵,v为初始向量,k为子空间的维数
%输出:lambda和U分别是A的k个特征值及相应的特征向量
tol=1.e-12;
v=v/norm(v); V(:,1)=v;
for i=1:k
    [V,T,beta]=Lanczos2(A,v,i);
    [Y,Mu]=eig(T);
    mu{i}=diag(Mu);
    U=V*Y;
    if beta*abs(Y(i,i))<tol
        k=i; break;
    end
end
lambda=mu{k};
```

例 7.12 令 $A \in \mathbb{R}^{500 \times 500}$ 是一个随机产生的对角矩阵, 其对角元服从正态分布. 对于 k 从 1 到 40, 利用算法 7.13 计算出 T_k 的特征值, 并标注在图 7.2 中. 图中第 k 列 “+” 表示 T_k 的所有特征值, 最后一列 “+” 表示 A 的精确特征值.

从这一数值例子的计算结果可以看出:

- (1) 两端特征值 (即最大和最小特征值) 收敛最快, 而内部特征值收敛较慢.
- (2) T_k 的第 i 个最大 (或最小) 特征值单调增加 (或减少) 地收敛到 A 的第 i 个最大 (或最小) 特征值.

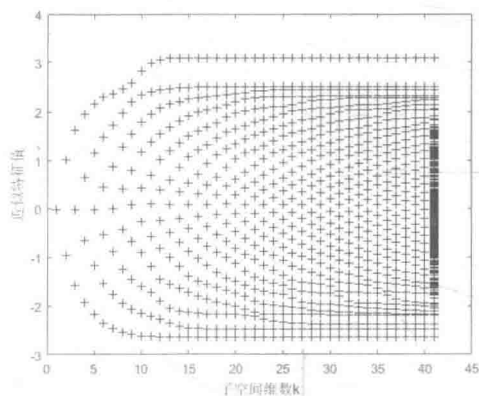


图 7.2 Lanczos 算法的收敛特性

2. 收敛性理论

下面对 Lanczos 方法及其收敛性进行一些理论上的分析. 为此, 先引进两个基本概念.

设 $x, y \in \mathbb{R}^n$ 是两个非零向量, 定义 x 与 y 之间的夹角为

$$\theta := \theta(x, y) = \arccos \frac{|x^T y|}{\|x\|_2 \|y\|_2}, \quad (7.101)$$

其几何解释如图 7.3 所示.

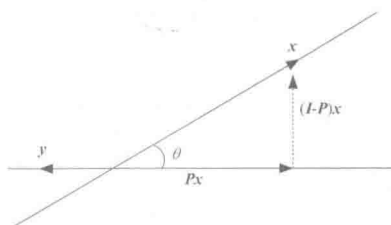


图 7.3 向量 x 与 y 之间的夹角

令

$$P = \frac{1}{\|y\|_2^2} y y^T,$$

即 P 是到 $\text{span}\{y\}$ 上的正交投影, 则有

$$\tan \theta = \frac{\|(I - P)x\|_2}{\|Px\|_2}, \quad \sin \theta = \frac{\|(I - P)x\|_2}{\|x\|_2}, \quad \cos \theta = \frac{\|Px\|_2}{\|x\|_2}.$$

设 $x \in \mathbb{R}^n$ 是一个非零向量, \mathcal{V} 是一个子空间, 则定义 x 与 \mathcal{V} 之间的夹角为

$$\theta(x, \mathcal{V}) = \min\{\theta(x, u) : 0 \neq u \in \mathcal{V}\}. \quad (7.102)$$

再假设 P 是 \mathcal{V} 上的正交投影, 则有

$$\tan \theta(x, \mathcal{V}) = \frac{\|(I - P)x\|_2}{\|Px\|_2}. \quad (7.103)$$

式 (7.103) 的几何解释如图 7.4 所示.

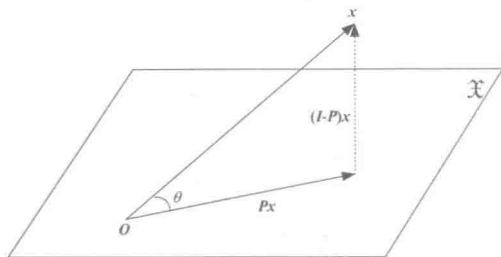


图 7.4 向量 x 与子空间 \mathcal{X} 之间的夹角

事实上, 可以证明 $\theta(x, \mathcal{V}) = \theta$, 其中 θ 如图 7.4 所示. 记 $x_1 = Px, x_2 = (I - P)x$, 则有

$$x = x_1 + x_2, \quad x_1^T x_2 = 0, \quad x_2 \perp \mathcal{V},$$

于是对任意的 $u \in \mathcal{V}$, 有

$$|x^T u| = |x_1^T u| \leq \|x_1\|_2 \|u\|_2,$$

从而

$$\frac{|x^T u|}{\|x\|_2 \|u\|_2} \leq \frac{\|x_1\|_2}{\|x\|_2} = \frac{|x_1^T x_1|}{\|x\|_2 \|x_1\|_2} = \frac{|x^T x_1|}{\|x\|_2 \|x_1\|_2}.$$

于是有

$$\theta(x, u) \geq \theta(x, x_1) = \theta.$$

注意到 u 的任意性, 便有 $\theta = \min_{u \in \mathcal{V}} \theta(x, u) = \theta(x, \mathcal{V})$.

为了叙述简洁, 先将下面所要使用的符号作一说明. 对给定的 $A^T = A \in \mathbb{R}^{n \times n}$ 和 $0 \neq v \in \mathbb{R}^n$, 简记

$$\mathcal{K}_k = \mathcal{K}_k(A, v),$$

并假定与其对应的长度为 k 的 Lanczos 分解为

$$AV_k = V_k T_k + \beta_k v_{k+1} e_k^T,$$

式中: $v_1 = v/\|v\|_2$. 再假定 A 和 T_k 的谱分解分别为

$$A = U \Lambda U^T \quad \text{和} \quad T_k = Y M Y^T,$$

式中:

$$\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n), \quad \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n,$$

$$M = \text{diag}(\mu_1, \mu_2, \dots, \mu_k), \quad \mu_1 \geq \mu_2 \geq \dots \geq \mu_k,$$

$$U = [u_1, u_2, \dots, u_n], \quad Y = [y_1, y_2, \dots, y_k].$$

如果没有特别说明, 在本节后续部分中所遇到的这些符号均为上面所述的含义.

首先给出一个简单的结果. 若 (λ, u) 是 A 的一个特征对, 而且如果 $u \in \mathcal{K}_k$, 则 λ 必是 T_k 的特征值. 事实上, $u \in \mathcal{K}_k$ 蕴涵着必存在 y 使得 $u = V_k y$, 这样利用 Lanczos 分解有

$$\lambda V_k y = \lambda u = Au = AV_k y = V_k T_k y + \beta_k v_{k+1} e_k^T y.$$

在上面的等式两边左乘 V_k^T , 并注意到 $V_k^T v_{k+1} = 0$, 有

$$\lambda y = T_k y,$$

即 λ 是 T_k 的特征值 (这里 $u \neq 0$ 蕴涵着 $y \neq 0$). 因此, 如果 $u \notin \mathcal{K}_k$, 但若 u 与 \mathcal{K}_k 很靠近的话, 可望 T_k 有特征值是 λ 的很好近似. 进一步, 可以证明如下结果.

定理 7.19 令 $Au = \lambda u$, $\|u\|_2 = 1$, 则必存在 T_k 的一个特征值 μ , 使得

$$|\mu - \lambda| \leq \|A\|_2 \tan \theta(u, \mathcal{K}_k), \quad (7.104)$$

式中: $\mathcal{K}_k = \mathcal{K}_k(A, v)$.

证明 令 $P = V_k V_k^T$ 是 \mathcal{K}_k 上的正交投影, 并记

$$u_1 = Pu, \quad u_2 = (I - P)u,$$

则有

$$\tan \theta(u, \mathcal{K}_k) = \frac{\|u_2\|_2}{\|u_1\|_2}, \quad (7.105)$$

而且

$$\begin{aligned} \lambda u_1 + \lambda u_2 &= \lambda u = Au = Au_1 + Au_2 \\ &= AV_k V_k^T u_1 + Au_2 \\ &= V_k T_k V_k^T u_1 + \beta_k v_{k+1} e_k^T V_k^T u_1 + Au_2. \end{aligned}$$

在上式两边左乘 V_k^T , 得

$$\lambda V_k^T u_1 = T_k V_k^T u_1 + V_k^T Au_2,$$

从而有

$$(\lambda I - T_k) V_k^T u_1 = V_k^T Au_2.$$

于是

$$\min_{\mu \in \lambda(T_k)} |\lambda - \mu| \cdot \|V_k^T u_1\|_2 = \min_{\mu \in \lambda(T_k)} |\lambda - \mu| \cdot \|Y^T V_k^T u_1\|_2$$

$$\leq \|(\lambda I - M)Y^T V_k^T u_1\|_2 \quad (7.106)$$

$$= \|(\lambda I - T_k)V_k^T u_1\|_2 = \|V_k^T A u_2\|_2$$

$$\leq \|V_k^T\|_2 \|A\|_2 \|u_2\|_2 = \|A\|_2 \|u_2\|_2.$$

此外, 由于 $u_1 \in \mathcal{K}_k$, 故存在 y 使得 $u_1 = V_k y$, 从而有

$$\|V_k^T u_1\|_2 = \|V_k^T V_k y\|_2 = \|y\|_2 = \|u_1\|_2, \quad (7.107)$$

将式 (7.105)、式 (7.106) 与式 (7.107) 相结合便得到所要证的不等式 (7.104). 证毕. \square

注 7.7 定理 7.19 是说, 如果对应于特征值 λ 的特征向量 u 与 Krylov 子空间 \mathcal{K}_k 之间的距离 $\tan \theta(u, \mathcal{K}_k)$ 很小, 则 T_k 就必有一个特征值 μ 和 λ 很靠近. 因此, 可望由 Lanczos 算法得到 λ 的一个很好的近似值.

下面的定理给出了第 i 个特征向量 u_i 与 \mathcal{K}_k 之间的距离 $\tan \theta(u_i, \mathcal{K}_k)$ 的一个可计算的上界估计.

定理 7.20 对给定的 $i (1 \leq i < k)$, 若

$$\lambda_{i-1} > \lambda_i, \quad \lambda_{i+1} > \lambda_n, \quad v^T u_i \neq 0,$$

则有

$$\tan \theta(u_i, \mathcal{K}_k) \leq \frac{\xi_i}{C_{k-i}(1 + 2\delta_i)} \tan \theta(u_i, v), \quad (7.108)$$

式中: $\mathcal{K}_k = \mathcal{K}_k(A, v)$; $\delta_i = \frac{\lambda_i - \lambda_{i+1}}{\lambda_{i+1} - \lambda_n}$; $C_{k-i}(t)$ 为 $k-i$ 次 Chebyshev 多项式; 而

$$\xi_i = \begin{cases} 1, & i = 1, \\ \prod_{s=1}^{i-1} \frac{\lambda_s - \lambda_n}{\lambda_s - \lambda_i}, & i > 1. \end{cases} \quad (7.109)$$

证明 记 $P_i = u_i u_i^T$, 即 P_i 是到子空间 $\text{span}\{u_i\}$ 上的正交投影, 因此, 有

$$\tan \theta(u_i, v) = \frac{\|(I - P_i)v\|_2}{\|P_i v\|_2} \quad (P_i v \neq 0). \quad (7.110)$$

此外, 任取 $x \in \mathcal{K}_k$, 则必存在一个 $p \in P_{k-1}$, 使得 $x = p(A)v$. 当然这里假定 $x \neq 0$. 注意到 $P_i A = A P_i$, 便有

$$P_i x = p(A)P_i v = p(\lambda_i)P_i v,$$

$$(I - P_i)x = p(A)(I - P_i)v = p(A)v_i,$$

式中: $v_i = (I - P_i)v$. 这样

$$\tan \theta(u_i, x) = \frac{\|(I - P_i)x\|_2}{\|P_i x\|_2} = \frac{\|p(A)v_i\|_2}{|p(\lambda_i)| \cdot \|P_i v\|_2}$$

$$\begin{aligned}
&= \left\| \frac{p(\mathbf{A})}{p(\lambda_i)} \frac{\mathbf{v}_i}{\|\mathbf{v}_i\|_2} \right\|_2 \frac{\|(I - \mathbf{P}_i)\mathbf{v}\|_2}{\|\mathbf{P}_i\mathbf{v}\|_2} \\
&= \|\widehat{p}(\mathbf{A})\widehat{\mathbf{v}}_i\|_2 \tan \theta(\mathbf{u}_i, \mathbf{v}),
\end{aligned} \tag{7.111}$$

式中: $\widehat{p}(t) = p(t)/p(\lambda_i)$; $\widehat{\mathbf{v}}_i = \mathbf{v}_i/\|\mathbf{v}_i\|_2$. 这里最后一个等式用到了式 (7.110).

现将 $\widehat{\mathbf{v}}_i$ 用 \mathbf{A} 的特征向量展开, 即令

$$\widehat{\mathbf{v}}_i = \alpha_1 \mathbf{u}_1 + \alpha_2 \mathbf{u}_2 + \cdots + \alpha_n \mathbf{u}_n.$$

而 $\widehat{\mathbf{v}}_i \perp \mathbf{u}_i$, 故必有 $\alpha_i = 0$. 此外, $\alpha_1^2 + \cdots + \alpha_n^2 = \|\widehat{\mathbf{v}}_i\|_2^2 = 1$. 这样, 有

$$\|\widehat{p}(\mathbf{A})\widehat{\mathbf{v}}_i\|_2^2 = \sum_{s=1}^n \alpha_s^2 \widehat{p}(\lambda_s)^2 \leq \max_{1 \leq s \leq n, s \neq i} |\widehat{p}(\lambda_s)|^2. \tag{7.112}$$

现在取一个特殊的多项式 $p_* \in \mathcal{P}_{k-1}$:

$$p_*(t) = \frac{(\lambda_1 - t)(\lambda_2 - t) \cdots (\lambda_{i-1} - t)}{(\lambda_1 - \lambda_i)(\lambda_2 - \lambda_i) \cdots (\lambda_{i-1} - \lambda_i)} v(t),$$

这里

$$v(t) = C_{k-i} \left(\frac{2t - (\lambda_{i+1} + \lambda_n)}{\lambda_{i+1} - \lambda_n} \right) / C_{k-i}(1 + 2\delta_i).$$

容易验证, 对于这个多项式 p_* , 有

$$p_*(\lambda_i) = 1, \quad p_*(\lambda_s) = 0, \quad s = 1, 2, \cdots, i-1.$$

从而有 $p_*(t) \equiv \widehat{p}_*(t)$, 而且

$$\begin{aligned}
\max_{1 \leq s \leq n, s \neq i} |p_*(\lambda_s)| &= \max_{i+1 \leq s \leq n} |p_*(\lambda_s)| \leq \xi_i \max_{i+1 \leq s \leq n} |v(\lambda_s)| \\
&\leq \frac{\xi_i}{C_{k-i}(1 + 2\delta_i)},
\end{aligned} \tag{7.113}$$

这里最后一个等式用到 Chebyshev 多项式的性质. 现将式 (7.111)、式 (7.112) 和式 (7.113) 相结合, 有

$$\begin{aligned}
\tan \theta(\mathbf{u}_i, \mathcal{K}_k) &= \min_{\mathbf{x} \in \mathcal{K}_k} \tan \theta(\mathbf{u}_i, \mathbf{x}) \\
&\leq \max_{1 \leq s \leq n, s \neq i} |p_*(\lambda_s)| \tan \theta(\mathbf{u}_i, \mathbf{v}) \\
&\leq \frac{\xi_i}{C_{k-i}(1 + 2\delta_i)} \tan \theta(\mathbf{u}_i, \mathbf{v}).
\end{aligned}$$

这正是要证的结论. 证毕. □

定理 7.21 (Kaniel-Saad 定理) 对于任意给定的 i ($1 \leq i < k$), 若

$$\lambda_i \neq \mu_s, \quad s = 1, 2, \cdots, i-1, \quad \lambda_{i+1} > \lambda_n,$$

而且 $v^T u_i \neq 0$, 则有

$$0 \leq \lambda_i - \mu_i \leq (\lambda_i - \lambda_n) \left[\frac{\zeta_i}{C_{k-i}(1+2\delta_i)} \tan \theta(u_i, v) \right]^2, \quad (7.114)$$

其中 δ_i 由定理 7.20 给出, 而

$$\zeta_i = \begin{cases} 1, & i = 1, \\ \max_{i+1 \leq \ell \leq n} \prod_{s=1}^{i-1} \frac{\mu_s - \lambda_\ell}{\mu_s - \lambda_i}, & i > 1. \end{cases} \quad (7.115)$$

注 7.8 不等式 (7.114) 右端可分为三部分: 第一部分为 $(\lambda_i - \lambda_n)\zeta_i^2$, 主要是由 λ_i 与 μ_1, \dots, μ_{i-1} 的分离程度确定; 第二部分为 $\tan^2 \theta(u_i, v)$, 是由初始向量 v 所确定, 反映了 v 中含有特征向量 u_i 的成分的大小; 第三部分为 $[C_{k-i}(1+2\delta_i)]^{-2}$, 这是关键的一部分, 反映了 Lanczos 方法的收敛性, 当 i 较小且 $0 < \delta_i < 0.1$ 时, 由 Chebyshev 多项式的性质可知, 此时随着 k 的增加, 这个值减少得非常快.

注 7.9 定理 7.21 实质上表明, 随着 Lanczos 迭代次数 k 的增加, T_k 的几个最大特征值 μ_1, \dots, μ_s (s 较小) 将非常快地收敛到 A 的前几个最大的特征值 $\lambda_1, \dots, \lambda_s$. 特别地, 对于 $i = 1$, 不等式 (7.114) 就变为

$$0 \leq \lambda_1 - \mu_1 \leq (\lambda_1 - \lambda_n) \tan^2 \theta(u_1, v) [C_{k-1}(1+2\delta_1)]^{-2},$$

由此更容易看出, 随着 k 的增加, μ_1 与 λ_1 之间的差距迅速地缩小.

注 7.10 对 $-A$ 应用定理 7.21 可知, T_k 的几个最小的特征值 $\mu_{k-s}, \mu_{k-s+1}, \dots, \mu_k$ 将随着 k 的增加很快地收敛到 A 的几个最小的特征值 $\lambda_{n-s}, \lambda_{n-s+1}, \dots, \lambda_n$. 特别地, 有

$$0 \leq \mu_k - \lambda_n \leq (\lambda_1 - \lambda_n) \tan^2 \theta(u_n, v) [C_{k-1}(1+2\tilde{\delta}_1)]^{-2},$$

式中: $\tilde{\delta}_1 = \frac{\lambda_n - 1 - \lambda_n}{\lambda_1 - \lambda_{n-1}}$. 随着 k 的增加, T_k 的最小特征值 μ_k 将很快地收敛到 A 的最小特征值 λ_n .

为了证明定理 7.21, 先将证明中需要的几个基本结果总结在下面的引理中.

引理 7.6 符号如前所述, 则有

(1) T_k 的每个特征值 μ_i 可表示为

$$\mu_i = \max_{0 \neq u \in \mathcal{U}_{k-i+1}} \frac{u^T A u}{u^T u}, \quad i = 1, 2, \dots, k, \quad (7.116)$$

式中:

$$\mathcal{U}_{k-i+1} = \text{span}\{V_k y_i, V_k y_{i+1}, \dots, V_k y_k\}. \quad (7.117)$$

(2) 由式 (7.117) 所定义的子空间 \mathcal{U}_{k-i+1} 可表示为

$$\mathcal{U}_{k-i+1} = \{p(A)v : p \in \mathcal{P}_{k-1}, p(\mu_s) = 0, s = 1, 2, \dots, i-1\}. \quad (7.118)$$

证明 (1) 由 T_k 的谱分解易证

$$\mu_i = \max_{0 \neq y \in \mathcal{Y}_{k-i+1}} \frac{y^T T_k y}{y^T y}, \quad (7.119)$$

式中: $\mathcal{Y}_{k-i+1} = \text{span}\{y_i, y_{i+1}, \dots, y_k\}$. 注意到

$$\frac{y^T T_k y}{y^T y} = \frac{y^T V_k^T A V_k y}{y^T V_k^T V_k y} = \frac{u^T A u}{u^T u}, \quad u = V_k y,$$

以及 $V_k \mathcal{Y}_{k-i+1} = \mathcal{U}_{k-i+1}$, 由式 (7.119) 可知式 (7.116) 成立.

(2) 由于

$$\mathcal{U}_{k-i+1} \subset \mathcal{R}(V_k) = \mathcal{K}_k = \{p(A)v : p \in \mathcal{P}_{k-1}\},$$

$$\mathcal{U}_{k-i+1}^\perp = \text{span}\{V_k y_1, V_k y_2, \dots, V_k y_{i-1}\},$$

故 $u \in \mathcal{U}_{k-i+1}$ 的充分必要条件是存在 $p \in \mathcal{P}_{k-1}$ 满足

$$v^T p(A) V_k y_s = 0, \quad s = 1, 2, \dots, i-1, \quad (7.120)$$

使得 $u = p(A)v$.

注意到 $v_{k+1} \perp \mathcal{K}_k$, 从 Lanczos 分解出发, 归纳地可以证明

$$v^T A^s V_k = v^T V_k T_k^s, \quad s = 0, 1, \dots, k-1.$$

从而, 对任意的 $p \in \mathcal{P}_{k-1}$, 有

$$v^T p(A) V_k = v^T V_k p(T_k). \quad (7.121)$$

上式两边右乘 y_s , 得

$$v^T p(A) V_k y_s = v^T V_k p(\mu_s) y_j = \|v\|_2 p(\mu_j) e_1^T y_s,$$

这里用到了 $v^T V_k = \|v\|_2 v_1^T V_k = \|v\|_2 e_1^T$. 再注意到 T_k 是不可约的对称三对角矩阵蕴含着 $e_1^T y_s \neq 0$, 便知式 (7.120) 成立的充分必要条件是

$$p(\mu_s) = 0, \quad s = 1, 2, \dots, i-1.$$

由此可知式 (7.118) 成立. 证毕. \square

有了前面的准备工作, 下面证明定理 7.21.

定理 7.21 的证明 注意到式 (7.94), 并且应用推论 7.2, 即知 $\lambda_i - \mu_i \geq 0$. 再应用引理 7.6 的结论 (1), 有

$$\begin{aligned} 0 \leq \lambda_i - \mu_i &= \lambda_i - \max_{0 \neq u \in \mathcal{U}_{k-i+1}} \frac{u^T A u}{u^T u} \\ &= \min_{0 \neq u \in \mathcal{U}_{k-i+1}} \frac{u^T (\lambda_i I - A) u}{u^T u}. \end{aligned} \quad (7.122)$$

现在定义 $\hat{p}(t) = g(t)h(t)$, 其中

$$g(t) = \begin{cases} 1, & i = 1, \\ \frac{(\mu_1 - t) \cdots (\mu_{i-1} - t)}{(\mu_1 - \lambda_i) \cdots (\mu_{i-1} - \lambda_i)}, & i > 1, \end{cases}$$

$$h(t) = C_{k-i} \left(\frac{2t - \lambda_{i+1} - \lambda_n}{\lambda_{i+1} - \lambda_n} \right) / C_{k-i}(1 + 2\delta_i).$$

容易验证, 这样定义的多项式 $\hat{p}(t)$ 满足 $\hat{p}(t) \in \mathcal{P}_{k-1}$, 且

$$\hat{p}(\lambda_i) = 1, \quad \hat{p}(\mu_s) = 0, \quad s = 1, 2, \cdots, i-1.$$

再利用引理 7.6 的结论 (2), 即有 $\hat{\mathbf{u}} = \hat{p}(\mathbf{A})\mathbf{v} \in \mathcal{U}_{k-i+1}$.

设 $\mathbf{v} = \gamma_1 \mathbf{u}_1 + \cdots + \gamma_n \mathbf{u}_n$. 直接计算有

$$\hat{\mathbf{u}}^T \hat{\mathbf{u}} = \mathbf{v}^T \hat{p}(\mathbf{A})^2 \mathbf{v} = \gamma_1^2 \hat{p}(\lambda_1)^2 + \cdots + \gamma_n^2 \hat{p}(\lambda_n)^2 \geq \gamma_i^2 \quad (7.123)$$

和

$$\begin{aligned} \hat{\mathbf{u}}^T (\lambda_i \mathbf{I} - \mathbf{A}) \hat{\mathbf{u}} &= \mathbf{v}^T \hat{p}(\mathbf{A}) (\lambda_i \mathbf{I} - \mathbf{A}) \hat{p}(\mathbf{A}) \mathbf{v} \\ &= \sum_{s=1}^n \gamma_s^2 \hat{p}(\lambda_s)^2 (\lambda_i - \lambda_s) \\ &= \sum_{s=1}^{i-1} \gamma_s^2 \hat{p}(\lambda_s)^2 (\lambda_i - \lambda_s) + \sum_{s=i+1}^n \gamma_s^2 \hat{p}(\lambda_s)^2 (\lambda_i - \lambda_s) \\ &\leq \sum_{s=i+1}^n \gamma_s^2 \hat{p}(\lambda_s)^2 (\lambda_i - \lambda_s) \\ &\leq (\lambda_i - \lambda_n) \sum_{s=i+1}^n \gamma_s^2 \hat{p}(\lambda_s)^2, \end{aligned} \quad (7.124)$$

其中不等式 (7.123) 用到了 $\hat{p}(\lambda_i) = 1$.

将式 (7.122)、式 (7.123) 和式 (7.124) 结合, 有

$$\begin{aligned} 0 \leq \lambda_i - \mu_i &\leq \frac{\hat{\mathbf{u}}^T (\lambda_i \mathbf{I} - \mathbf{A}) \hat{\mathbf{u}}}{\hat{\mathbf{u}}^T \hat{\mathbf{u}}} \\ &\leq \frac{\lambda_i - \lambda_n}{\gamma_i^2} \sum_{s=i+1}^n \gamma_s^2 \hat{p}(\lambda_s)^2 \\ &\leq \frac{\lambda_i - \lambda_n}{\gamma_i^2} \max_{i+1 \leq s \leq n} \hat{p}(\lambda_s)^2 \sum_{s=i+1}^n \gamma_s^2. \end{aligned} \quad (7.125)$$

此外, 由 $\hat{p}(t)$ 的定义, 有

$$\max_{i+1 \leq s \leq n} |\hat{p}(\lambda_s)| \leq \max_{i+1 \leq s \leq n} |g(\lambda_s)| \cdot \max_{i+1 \leq s \leq n} |h(\lambda_s)|$$

$$\leq \frac{\zeta_i}{C_{k-i}(1+2\delta_i)}. \quad (7.126)$$

将式 (7.126) 代入式 (7.125), 并注意到

$$\begin{aligned} \frac{1}{\gamma_i^2} \sum_{s=i+1}^n \gamma_s^2 &\leq \frac{1}{\gamma_i^2} \left(\sum_{s=1}^{i-1} \gamma_s^2 + \sum_{s=i+1}^n \gamma_s^2 \right) \\ &= \frac{\|(I - P_i)v\|_2^2}{\|P_i v\|_2^2} = \tan^2 \theta(u_i, v), \end{aligned}$$

便知式 (7.114) 成立. 证毕. \square

3. 重开始 Lanczos 方法

从理论上讲, 由 Lanczos 算法产生的矩阵 V_k 的列向量是相互正交的. 然而, 当其在计算机上运行时, 产生的 V_k 很快就失去了其正交性. 正交性的损失会导致 Ritz 值的收敛发生紊乱现象: 会有多个 Ritz 值收敛到同一个特征值, 或者有的特征值迟迟没有逼近它的 Ritz 值出现. 如果并不关心特征值的重数, 又能忍受长时间的等待, 则正交性的损失并不影响所要达到的目的. 但如果对特征值精确的重数很感兴趣, 或者对计算时间要求很高, 则就必须去弥补 V_k 的正交性损失. 目前已有各种各样的重正交化技术来达到这一目标, 最直接的方法就是采用算法 2.14 后所述的完全重正交化技术. 这样做的结果是 V_k 的正交性可以达到机器精度, 但为此付出的代价是:

- (1) 存储量从 $O(n)$ 增加到 $O(kn)$.
- (2) 计算量从 $O(kn)$ 增加到 $O(k^2n)$.

由于计算机存储空间的限制, k 不能无限制地增加, 这样就会出现 k 已经达到了最大的可能, 但 T_k 中还没有足够的信息来解决希望解决的问题.

解决这一问题的方法主要有两种: 一种是放弃完全重正交化, 而采用有选择性的重正交化技术, 从而减少存储量的要求; 另一种是使用重开始技术.

假定已经得到了一个长度为 k 的 Lanczos 分解

$$AV_k = V_k T_k + \beta_k v_{k+1} e_k^T, \quad (7.127)$$

而希望求的特征值是 \hat{m} 个 (如 \hat{m} 个最大特征值, 或者 \hat{m} 个最小特征值). 现在取一个正整数 m 略大于 \hat{m} , 如 $m = \hat{m} + 2$. k 常取作 $2m$ 或 $5m$ 等.

下面介绍最近由 Stewart 给出的 Krylov-Schur 重开始方法, 其基本步骤如下.

第 1 步, 计算 T_k 的谱分解 $T_k = YMY^T$, 其中

$$Y = [y_1, y_2, \dots, y_k], \quad M = \text{diag}(\mu_1, \mu_2, \dots, \mu_k),$$

然后适当地调整特征值 μ_i 的次序, 使得 μ_{m+1}, \dots, μ_k 是与要计算的特征值差距较大者, 希望将其剔除. 例如, 要计算 A 的几个最大特征值, 则将 μ_i 排序为

$$\mu_1 \geq \mu_2 \geq \dots \geq \mu_m \geq \dots \geq \mu_k$$

即可.

第 2 步, 将 T_k 的谱分解代入式 (7.127), 并右乘 Y , 得

$$A(V_k Y) = (V_k Y)M + \beta_k v_{k+1} e_k^T Y.$$

比较上式两边的前 m 列, 得

$$A\tilde{V}_m = \tilde{V}_m M_1 + \beta_k v_{k+1} z_m^T, \quad (7.128)$$

式中:

$$\tilde{V}_m = V_k Y(:, 1:m), \quad z_m^T = Y(k, 1:m), \quad M_1 = \text{diag}(\mu_1, \dots, \mu_m).$$

第 3 步, 计算一个正交矩阵 $Q \in \mathbb{R}^{m \times m}$, 使得

$$QM_1Q^T = \hat{T}_m, \quad Qz_m = \alpha e_m, \quad (7.129)$$

式中: \hat{T}_m 为对称三对角矩阵; $\alpha = \|z_m\|_2$. 则在式 (7.128) 两边右乘 Q^T , 可得如下长度为 m 的 Lanczos 分解:

$$A\hat{V}_m = \hat{V}_m \hat{T}_m + \hat{\beta}_m v_{m+1} e_m^T, \quad (7.130)$$

式中:

$$\hat{V}_m = \tilde{V}_m Q^T, \quad \hat{\beta}_m = \alpha \beta_k, \quad v_{m+1} = v_{k+1}.$$

然后从式 (7.130) 出发, 再应用重正交化的 Lanczos 迭代, 使其扩展为一个长度为 k 的 Lanczos 分解. 如果仍不满足要求, 则可再用上面介绍的方法将其收缩为一个长度为 m 的 Lanczos 分解. 如此反复进行, 最终会得到一个 \hat{T}_m , 它的特征值含有所需要的特征值的很好的近似.

现在的问题是如何实现式 (7.129). 这一计算任务的具体计算过程可从下面的示例中明白. 例如 $m = 5$, 且已经将 M_1 和 z_m 约化为如下形式, 即

$$M_1 \rightarrow \tilde{M}_1 = \begin{bmatrix} \alpha_1 & \beta_1 & & & \\ & \beta_1 & \alpha_2 & \beta_2 & \\ & & \beta_2 & \alpha_3 & \\ & & & & \alpha_4 \\ & & & & & \alpha_5 \end{bmatrix}, \quad z_m \rightarrow \tilde{z}_m = \begin{bmatrix} 0 \\ 0 \\ \tilde{z}_3 \\ z_4 \\ z_5 \end{bmatrix}.$$

下一步, 先确定一个 Givens 变换 $G_3 = G(3, 4; c_3, s_3)$, 使得 \tilde{z}_m 的第 3 个分量为零, 然后计算 $\hat{z}_m = G_3 \tilde{z}_m$ 和 $\hat{M}_1 = G_3 \tilde{M}_1 G_3^T$, 则 \hat{z}_m 和 \hat{M}_1 有如下形式, 即

$$\hat{M}_1 = \begin{bmatrix} \alpha_1 & \beta_1 & & & \\ \beta_1 & \alpha_2 & \beta_2 & \gamma & \\ & \beta_2 & \alpha_3 & \beta_3 & \\ & & \gamma & \beta_3 & \alpha_4 \\ & & & & & \alpha_5 \end{bmatrix}, \quad \hat{z}_m = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \tilde{z}_4 \\ z_5 \end{bmatrix}.$$

在 \widehat{M}_1 的 (2,4) 和 (4,2) 位置上出现了两个不希望有的可能非零元 “ γ ”. 然后, 就可以采用类似于隐式对称 QR 方法中使用的技术, 计算 (2,3) 平面和 (1,2) 平面的两个 Givens 变换

$$G_2 = G(2, 3; c_2, s_2) \text{ 和 } G_1 = G(1, 2; c_1, s_1),$$

将这两个不需要的元素消去, 即

$$\widehat{M}_1 \xrightarrow{G_2} \begin{bmatrix} \alpha_1 & \beta_1 & \gamma & & \\ \beta_1 & \alpha_2 & \beta_2 & & \\ \gamma & \beta_2 & \alpha_3 & \beta_3 & \\ & & \beta_3 & \alpha_4 & \\ & & & & \alpha_5 \end{bmatrix} \xrightarrow{G_1} \begin{bmatrix} \alpha_1 & \beta_1 & & & \\ \beta_1 & \alpha_2 & \beta_2 & & \\ & \beta_2 & \alpha_3 & \beta_3 & \\ & & \beta_3 & \alpha_4 & \\ & & & & \alpha_5 \end{bmatrix}.$$

再注意到 $G_1 G_2 \widehat{z}_m = \widehat{z}_m = \widehat{z}_m$, 就完成了步约化, 即在 z_m 上又多引进了一个零元. 具体的算法可总结如下.

算法 7.14 给定 m 阶对角矩阵 $M = \text{diag}(\mu_1, \dots, \mu_m)$ 和 m 维向量 z . 本算法计算正交矩阵 Q , 对称三对角矩阵 T 和实数 α , 使得 $Qz = \tau e_m$, $QM Q^T = T$.

function $[Q, T, \tau] = \text{TriReduce}(M, z)$

$\alpha_i = \mu_i, i = 1, 2, \dots, m;$

$\beta_i = 0, i = 1, 2, \dots, m-1;$

$Q = I_m;$

for $i = 1 : m-1$

确定 $c = \cos \theta$ 和 $s = \sin \theta$, 使得

$$\begin{bmatrix} c & s \\ -s & c \end{bmatrix} \begin{bmatrix} z_i \\ z_{i+1} \end{bmatrix} = \begin{bmatrix} 0 \\ \delta \end{bmatrix}; \quad \begin{bmatrix} z_i \\ z_{i+1} \end{bmatrix} := \begin{bmatrix} 0 \\ \delta \end{bmatrix};$$

$$\begin{bmatrix} \alpha_i & \beta_i \\ \beta_i & \alpha_{i+1} \end{bmatrix} := \begin{bmatrix} c & s \\ -s & c \end{bmatrix} \begin{bmatrix} \alpha_i & \beta_i \\ \beta_i & \alpha_{i+1} \end{bmatrix} \begin{bmatrix} c & s \\ -s & c \end{bmatrix}^T;$$

$\gamma := -s\beta_{i-1}; \beta_{i-1} := c\beta_{i-1};$

$Q := G(i, i+1; c, s)Q;$

if $i > 1$

for $r = i : -1 : 2$

确定 $c = \cos \theta$ 和 $s = \sin \theta$, 使得

$$\begin{bmatrix} c & s \\ -s & c \end{bmatrix} \begin{bmatrix} \gamma \\ \beta_r \end{bmatrix} = \begin{bmatrix} 0 \\ \sigma \end{bmatrix}; \quad \beta_r := \sigma;$$

$$\begin{bmatrix} \alpha_{r-1} & \beta_{r-1} \\ \beta_{r-1} & \alpha_r \end{bmatrix} := \begin{bmatrix} c & s \\ -s & c \end{bmatrix} \begin{bmatrix} \alpha_{r-1} & \beta_{r-1} \\ \beta_{r-1} & \alpha_r \end{bmatrix} \begin{bmatrix} c & s \\ -s & c \end{bmatrix}^T;$$

if $r > 2$

$$\gamma := -s\beta_{r-2}; \quad \beta_{r-2} := c\beta_{r-2};$$

end

$$Q := G(r-1, r; c, s)Q;$$

end

end

end

$$\tau = z_m;$$

注 7.11 需要注意的是, 这里对于对称三对角矩阵 T , 只存储其对角元和次对角元: $\alpha_1, \alpha_2, \dots, \alpha_m; \beta_1, \beta_2, \dots, \beta_{m-1}$.

算法 7.14 的 MATLAB 程序如下:

```
function [Q,alpha,beta,tau]=trireduce(M,z)
%输入:M为m阶对角矩阵,z为m维向量
%输出:Q为正交矩阵,alpha为对称三对角矩阵T的对角元,
%      beta为对称三对角矩阵T的次对角元,使得Qz=tau*e_m,QMQ'=T
m=size(M,1);
for i=1:m
    alpha(i)=M(i,i);
end
beta=zeros(m-1,1); Q=eye(m);
for i=1:m-1
    delta=sqrt(z(i)^2+z(i+1)^2);
    c=z(i+1)/delta; s=-z(i)/delta;
    z(i+1)=delta; z(i)=0;
    a=[c,s;-s,c]*[alpha(i),beta(i);
    beta(i),alpha(i+1)]*[c,s;-s,c]';
    alpha(i)=a(1,1); beta(i)=a(1,2); alpha(i+1)=a(2,2);
    Q(i:i+1,:)=[c,s;-s,c]*Q(i:i+1,:);
    if i>1
        gamma=-s*beta(i-1); beta(i-1)=c*beta(i-1);
        for r=i:-1:2
            sigma=sqrt(gamma^2+beta(r)^2);
            c=beta(r)/sigma; s=-gamma/sigma;
            a=[c,s;-s,c]*[alpha(r-1),beta(r-1);
            beta(r-1),alpha(r)]*[c,s;-s,c]';
```

```

        alpha(r-1)=a(1,1); beta(r-1)=a(1,2);
        alpha(r)=a(2,2); beta(r)=sigma;
        if r>2
            gamma=-s*beta(r-2); beta(r-2)=c*beta(r-2);
        end
        Q(r-1:r,:)= [c,s;-s,c]*Q(r-1:r,:);
    end
end
end
tau=z(m);

```

在上面的讨论中, 已经提到了用 Lanczos 过程将长度为 m 的 Lanczos 分解扩展到长度为 k 的 Lanczos 分解, 下面算法给出了一种实现方法.

算法 7.15 给定对称矩阵 A 的长度为 m 的 Lanczos 分解

$$AV_m = V_m T_m + \beta_m v_{m+1} e_m^T = V_{m+1} \hat{T}_m,$$

本算法将其扩展为一个长度为 k 的 Lanczos 分解.

function $[V_{k+1}, \hat{T}_k] = \text{ExpandLan}(A, V_{m+1}, \hat{T}_m, m, k)$

for $i = m+1 : k$

$u = Av_i - \beta_{i-1}v_{i-1}; \alpha_i = u^T v_i;$

$u = u - \alpha_i v_i;$

$u = u - (v_{i-1}^T u) v_{i-1};$ (局部重正交化)

$u = u - (v_i^T u) v_i;$

$u = u - V_i(V_i^T u);$ (完全重正交化)

$\beta_i = \|u\|_2;$

if $\beta_i = 0$

stop; (已经得到了 A 的一个不变子空间)

else

$v_{i+1} = u/\beta_i;$

$V_{i+1} = [V_i, v_{i+1}];$

end

end

算法 7.15 的 MATLAB 程序如下:

function $[V, \alpha, \beta] = \text{expand_lanczos}(A, V, \alpha, \beta, m, k)$

%给定对称矩阵A的长度为m的Lanczos分解,

%本算法将其扩展为一个长度为k的Lanczos分解

for $i=m+1:k$

$u=A*V(:,i)-\beta(i-1)*V(:,i-1);$

$\alpha(i)=u'*V(:,i); u=u-\alpha(i)*V(:,i);$

```

u=u-(u'*V(:,i-1))*V(:,i-1); %局部重正交化
u=u-(u'*V(:,i))*V(:,i);
u=u-V(:,1:i)*(V(:,1:i)'\u); %完全重正交化
beta(i)=norm(u);
if (beta(i)==0)
    return;
else
    V(:,i+1)=u/beta(i);
    V=[V, V(:,i+1)];
end
end
end

```

综合上面的讨论, 就得到了如下的隐式重开始完全重正交化 Lanczos 方法.

算法 7.16 (隐式重开始 Lanczos 方法) 给定一个大型稀疏的对称矩阵 $A \in \mathbb{R}^{n \times n}$. 本算法计算 A 的 \hat{m} 个两端特征对 ($\hat{m} \ll n$).

步 1, 初始化

- (1) $m = \hat{m} + 2$; $k = 2m$ (或者 $5m$, 或者取存储所能允许的最大整数).
- (2) 选取一个适当的非零向量 $v \in \mathbb{R}^n$ (一般随机选取).
- (3) 计算

$$\begin{aligned}
 v_1 &= v / \|v\|_2; \quad u = Av_1; \quad \alpha_1 = u^T v_1; \\
 u &= u - \alpha_1 v_1; \quad u = u - (u^T v_1) v_1; \\
 \beta_1 &= \|u\|_2; \quad v_2 = u / \beta_1; \quad V_2 = [v_1, v_2].
 \end{aligned}$$

- (4) $[V_{m+1}, \hat{T}_m] = \text{ExpandLan}(A, V_2, \hat{T}_1, 1, m)$.

步 2, 迭代

- (1) $[V_{k+1}, \hat{T}_k] = \text{ExpandLan}(A, V_{m+1}, \hat{T}_m, m, k)$.
- (2) 用对称 QR 方法 (或其他方法) 计算 T_k 的谱分解:

$$T_k = YMY^T, \quad M = \text{diag}(\mu_1, \dots, \mu_k), \quad Y^T Y = I_k.$$

(3) 若 T_k 的特征值中已含有所求特征值的很好近似, 则输出有关消息, 结束; 否则适当调整 μ_i 的顺序, 使得 T_k 的两端特征值位于前面, 而后 $k-m$ 个特征值为需要剔除的对象.

- (4) $[Q, T_m, \alpha] = \text{TriReduce}(M_1, z)$, 其中

$$M_1 = \text{diag}(\mu_1, \dots, \mu_k), \quad z^T = Y(k, 1:m).$$

- (5) $V_{m+1} = [V_k Y(1:k, 1:m) Q^T, v_{k+1}]$; $\beta_m = \alpha \beta_k$. 然后, 转步 2.

注 7.12 假如希望计算的是 A 之位于某一实数 μ 附近的几个特征对, 则需要应用上一节最后所介绍的位移求逆技术, 即应用算法 7.16 到矩阵 $(A - \mu I)^{-1}$ 上.

算法 7.16 的 MATLAB 程序如下:

```
function [M,V,iter]=irest_Lanczos(A,m)
%给定一个大型稀疏的对称矩阵A,
%本算法计算A的m个两端特征对(m<<n).
m=m+2; k=2*m; V=[]; tol=1.0e-30; iter=0;
v=rand(size(A,1),1); v=v/norm(v); V=[V,v];
u=A*v; alpha(1)=u'*v;
u=u-alpha(1)*v; u=u-(u'*v)*v;
beta(1)=norm(u); v=u/beta(1); V=[V,v];
[V,alpha,beta]=expand_lanczos(A,V,alpha,beta,1,m);
while(iter<1000)
    iter=iter+1;
    [V,alpha,beta]=expand_lanczos(A,V,alpha,beta,m,k);
    betak=beta(end); %将beta的最后一个分量存起来
    vk1=V(:,end); %将V的最后一列存起来
    T=diag(alpha)+diag(beta(1:end-1),1)+diag(beta(1:end-1),-1);
    [Y,M]=eig(T); %M的对角元按升序按列
    [M,I]=sort(diag(M),'descend'); %对M的对角元按降序按列
    M=diag(M); Y=Y(:,I); %第k列到了第1列
    V=V(:,1:k)*Y(:,1:m); z=Y(k,1:m)'; M=M(1:m,1:m);
    if abs(beta(end))*abs(Y(k,1))<tol
        break;
    end
    [V,alpha,beta,theta]=trireduce2(M,z);
    V=V*V'; V=[V,vk1];
    beta(m)=theta*betak;
end
```

例 7.13 用算法 7.16 计算矩阵

$$A = \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & 2 & 2 & \cdots & 2 \\ 1 & 2 & 3 & \cdots & 3 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 2 & 3 & \cdots & n \end{bmatrix} \quad (7.131)$$

的 m 个模最大的特征值. 取 $n = 2000$, $m = 5$.

解 利用算法 7.16, 编写 MATLAB 脚本文件 ex713.m,

```
%例7.13-ex713.m
clear all
n=2000; m=5; A=zeros(n,n);
for i=1:n
    for j=1:i
        A(i,j)=j; A(j,i)=A(i,j);
    end
end
ev=eig(A); ev=sort(ev,'descend'); lam=ev(1:m);
[M,V,iter]=irest_Lanczos(A,m); iter
mu=diag(M(1:m,1:m))
err1=norm(mu-lam)
err2=norm(A*V(:,1)-M(1,1)*V(:,1))
```

然后在 MATLAB 命令窗口运行之, 即得所要求的结果.

7.6.3 Arnoldi 方法

Arnoldi 方法是最经典的一类 Krylov 子空间法, 它主要用于计算一个大型稀疏非对称矩阵的少数几个特征值和对应的特征向量. 随着重开始技术的不断完善, Arnoldi 方法的应用范围变得越来越广. 本节将阐述这些重开始技术以及与此相结合所产生的重开始 Arnoldi 方法.

1. 经典 Arnoldi 方法

设 $A \in \mathbb{R}^{n \times n}$ 已经给定, 希望计算的是 A 在某一指定范围内的少数几个特征值及对应的特征向量, 其中 A 是大型的稀疏矩阵. 将算法 2.12 与 Rayleigh-Ritz 方法相结合就可得到完成这一任务的 Arnoldi 算法.

算法 7.17 (经典 Arnoldi 方法) 给定一个 n 阶实矩阵 A . 本算法计算 A 的少数几个特征对.

步 1, 选择初始向量 v , 并令 $v_1 = v/\|v\|_2$.

步 2, 利用算法 2.12 产生一个长度为 k 的 Arnoldi 分解, 即

$$AV_k = V_k H_k + \beta_k v_{k+1} e_k^T.$$

步 3, 计算 H_k 的特征值, 并选择其中若干个满足要求的记为 μ_1, \dots, μ_l .

步 4, 计算 μ_i 所对应的特征向量 y_i , 并形成 Ritz 向量 $u_i = V_k y_i, i = 1, 2, \dots, l$.

步 5, 如果不满足要求, 增加 k , 再返回到步 2, 重复以上各步.

注 7.13 算法 7.17 只涉及矩阵乘向量 $y = Ax$, 因此可充分利用 A 的稀疏性 A 所具有的特殊结构.

注 7.14 设 $H_k y = \mu y$, 在 Arnoldi 分解的两边右乘 y , 得

$$A V_k y = \mu V_k y + \beta_k v_{k+1} e_k^T y,$$

于是有

$$\|A(V_k y) - \mu(V_k y)\|_2 = |\beta_k| \cdot |e_k^T y|. \quad (7.132)$$

由此可知, 可以由 $|\beta_k| \cdot |e_k^T y|$ 的大小来决定 $(\mu, V_k y)$ 是否可以作为 A 的近似特征对.

为了对算法 7.17 的数值特性有个直观的了解, 下面看一个简单的数值例子.

例 7.14 令

$$A = P^{-1} \text{diag}(-24, -23, \dots, 24, 25) P \in \mathbb{R}^{50 \times 50},$$

式中: P 为随机产生的可逆矩阵.

对于 k 从 1 到 30, 应用算法 7.17 于给定的矩阵 A 上, 并且将计算结果标注在了图 7.5 中, 其中横坐标表示迭代步数 k , 第 k 列的点表示 H_k 的所有实特征值, 最后一列的点表示 A 的精确特征值.

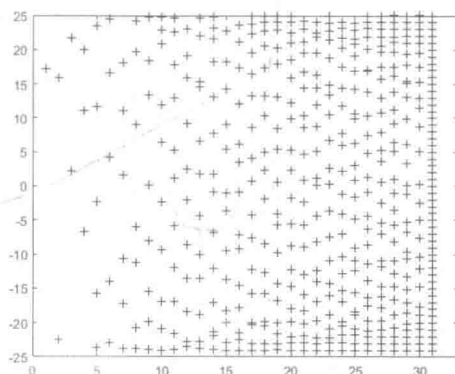


图 7.5 Arnoldi 算法的收敛特性

请注意观察该例所展示出的算法 7.17 的收敛特性: $k=10$ 时, H_{10} 就有两个最大特征值收敛到了 A 的两个最大特征值; 而到了 $k=21$ 时, H_{21} 就有 3 个最小特征值和 2 个最大特征值分别收敛到了 A 的 3 个最小特征值和 2 个最大特征值; 而位于区间内部的有些特征值, 直到 $k=30$ 还没有 Ritz 值收敛到它们.

虽然例 7.14 仅仅是一个可对角化的特殊矩阵, 但其所展示出的算法 7.17 的收敛特性对一般情形也是成立的, 即外围特征值 (exterior eigenvalues) 最先收敛, 而内部特征值收敛特别缓慢. 事实上, 这一结果也可以从理论上解释如下.

为了叙述简洁起见, 假设 A 是可对角化的. 设 A 的特征对为

$$(\lambda_j, x_j), \quad j = 1, 2, \dots, n,$$

则 A 可对角化的假定蕴含着 x_1, x_2, \dots, x_n 是线性无关的. 现将初始向量 v 按特征向量展开:

$$v = \gamma_1 x_1 + \gamma_2 x_2 + \dots + \gamma_n x_n.$$

根据定理 2.10 的结论 (5), $u \in \mathcal{K}_k(A, v)$ 的充分必要条件是存在 $p \in \mathcal{P}_{k-1}$, 使得 $u = p(A)v$, 从而有

$$u = \gamma_1 p(\lambda_1) x_1 + \gamma_2 p(\lambda_2) x_2 + \dots + \gamma_n p(\lambda_n) x_n.$$

由此可知, 欲使 u 是某一特征向量 x_j 的很好近似, 则必须有 $|p(\lambda_j)|$ 相对要比其他的 $|p(\lambda_i)|$ 大得多. 而由复变函数的最大模原理可知, 这只有 λ_j 位于谱集的边缘才能达到. 这也就是 Arnoldi 方法为什么先收敛到外围特征值的缘故.

前面已经列举了 Arnoldi 算法的优点, 但它也有致命的不足之处. 这一算法随着 k 的增加其存储空间需求越来越大, 因此 k 的大小受到了计算机存储量的限制. 这样一来就会出现这样的问题: 希望计算的特征值还没有得到, 而 k 就不能再增大了. 解决这一问题的方法就是现在已经逐渐成熟的重开始技术.

2. 隐式重开始 Arnoldi 方法

假设希望计算 A 的 m 个模最大的特征值 (其他情况可通过注 7.15 所述的“位移求逆技术”归结为这种情形). 再假定已得到一个长度为 k 的 Arnoldi 分解

$$AV_k = V_k H_k + \beta_k v_{k+1} e_k^T, \quad \beta_k = h_{k+1,k}. \quad (7.133)$$

由于计算机存储空间的限制, k 已经不能再增加了. 但由 H_k 所提供的 Ritz 值还不满足精度要求. 面临这样的处境, 下一步没有别的选择, 只有重新开始, 即再选择一个新的初始向量 \hat{v} , 重新再来计算新的 Arnoldi 分解. 但是并不想把前面所做的一切全部扔掉, 最好能够对重新选择 \hat{v} 提供一些有用的信息, 使新选择的 \hat{v} 比 v 更好一些. 这正是发展重开始技术的最原始想法.

目前存在两种隐式重开始 Arnoldi 方法: 一种是 Sorensen (索仑森) 提出的基于“过滤多项式”技术的隐式重开始方法; 另一种是 Stewart (史都瓦) 提出的基于实 Schur 分解的所谓 Krylov-Schur 重开始方法. 在此着重介绍后者, 其基本步骤有如下 3 步:

第 1 步, 计算 H_k 的实 Schur 分解: $H_k = U_1 R U_1^T$, 其中 $U_1 \in \mathbb{R}^{k \times k}$ 是正交矩阵, $R \in \mathbb{R}^{k \times k}$ 是拟上三角矩阵 (即 R 是分块上三角矩阵, 其对角块是 1×1 或 2×2). 这一步可由著名的 QR 方法实现.

第 2 步, 重排 R 的对角块, 即计算一个正交矩阵 $U_2 \in \mathbb{R}^{k \times k}$, 使得

$$U_2^T R U_2 = \begin{bmatrix} R_1 & * \\ O & R_2 \end{bmatrix},$$

式中: R_1 和 R_2 均为实 Schur 标准形, 而且

$$\lambda(R_1) = \{\mu_1, \dots, \mu_m\}, \quad \lambda(R_2) = \{\mu_{m+1}, \dots, \mu_k\}.$$

第 1 步和第 2 步计算完成之后, 有

$$\mathbf{H}_k = \mathbf{U}_1 \mathbf{U}_2 \begin{bmatrix} \mathbf{R}_1 & * \\ \mathbf{O} & \mathbf{R}_2 \end{bmatrix} (\mathbf{U}_1 \mathbf{U}_2)^T.$$

将其代入式 (7.133) 并且右乘 $\mathbf{U}_1 \mathbf{U}_2$, 得

$$\mathbf{A}(\mathbf{V}_k \mathbf{U}_1 \mathbf{U}_2) = \mathbf{V}_k \mathbf{U}_1 \mathbf{U}_2 \begin{bmatrix} \mathbf{R}_1 & * \\ \mathbf{O} & \mathbf{R}_2 \end{bmatrix} + \beta_k \mathbf{v}_{k+1} \mathbf{e}_k^T (\mathbf{U}_1 \mathbf{U}_2).$$

比较上式两边的前 m 列, 得

$$\mathbf{A} \bar{\mathbf{V}}_m = \bar{\mathbf{V}}_m \mathbf{R}_1 + \beta_k \mathbf{v}_{k+1} \mathbf{z}_m^T, \quad (7.134)$$

式中: $\bar{\mathbf{V}}_m$ 为 $\mathbf{V}_k \mathbf{U}_1 \mathbf{U}_2$ 的前 m 列构成的矩阵; \mathbf{z}_m^T 为 $\mathbf{U}_1 \mathbf{U}_2$ 的最后一行的前 m 个元素构成的向量.

第 3 步, 计算一个正交矩阵 $\mathbf{Q} \in \mathbb{R}^{m \times m}$, 使得

$$\mathbf{Q}^T \mathbf{R}_1 \mathbf{Q} = \widehat{\mathbf{H}}_m, \quad \mathbf{z}_m^T \mathbf{Q} = \alpha \mathbf{e}_m^T, \quad (7.135)$$

式中: $\widehat{\mathbf{H}}_m$ 为上 Hessenberg 矩阵.

事实上, 式 (7.135) 可按如下方式实现:

(1) 首先计算一个 Householder 变换 \mathbf{P}_1 , 使得

$$\mathbf{z}_m^T \mathbf{P}_1 = \alpha \mathbf{e}_m^T.$$

(2) 然后依次计算 $m-2$ 个 Householder 变换

$$\mathbf{P}_i \in \mathbb{R}^{(m-i-1) \times (m-i+1)}, \quad i = 2, \dots, m-1,$$

使得

$$\mathbf{Q}_{m-1}^T \cdots \mathbf{Q}_2^T \mathbf{P}_1^T \mathbf{R}_1 \mathbf{P}_1 \mathbf{Q}_2 \cdots \mathbf{Q}_{m-1} = \widehat{\mathbf{H}}_m$$

为上 Hessenberg 矩阵, 其中 $\mathbf{Q}_i = \text{diag}(\mathbf{P}_i, \mathbf{I}_{i-1})$.

其实, 这一约化过程本质上就是在 7.4.1 节中介绍的约化一个矩阵为上 Hessenberg 形的方法. 只是此处是由矩阵的行向量来确定约化所需的 Householder 变换 \mathbf{P}_i , 而且是从最后一行开始约化的; 而那里是由其列向量来确定, 而且是从第 1 列开始约化的. 显然, 这样得到的 \mathbf{Q}_i 满足 $\mathbf{e}_m^T (\mathbf{Q}_2 \cdots \mathbf{Q}_{m-1}) = \mathbf{e}_m^T$, 故令 $\mathbf{Q} = \mathbf{P}_1 \mathbf{Q}_2 \cdots \mathbf{Q}_{m-1}$, 则该正交矩阵就是式 (7.135) 中所需的正交矩阵.

一旦式 (7.135) 已经得到, 则在式 (7.134) 两边右乘 \mathbf{Q} , 得

$$\mathbf{A} \widehat{\mathbf{V}}_m = \widehat{\mathbf{V}}_m \widehat{\mathbf{H}}_m + \widehat{\beta}_m \widehat{\mathbf{v}}_{m+1} \mathbf{e}_m^T, \quad (7.136)$$

式中: $\widehat{\mathbf{V}}_m = \bar{\mathbf{V}}_m \mathbf{Q}$, $\widehat{\mathbf{v}}_{m+1} = \mathbf{v}_{k+1}$, $\widehat{\beta}_m = \alpha \beta_k$.

这就是 Stewart 的 Krylov-Schur 重新开始方法最终计算得到的结果: 将一个长度为 k 的 Arnoldi 分解收缩为一个长度为 m 的 Arnoldi 分解.

下面考虑 Krylov-Schur 重开始方法的收敛性问题. 为了符号简单, 将式 (7.134) 改写为

$$AV_m = V_m R + \beta_k v_{k+1} z^T. \quad (7.137)$$

设

$$z = (0, \dots, 0, z_{\ell+1}, \dots, z_m)^T. \quad (7.138)$$

实际计算时, 要给定一个收敛准则, 可将 z 之小到一定程度的分量均置零. 现假定对此 ℓ , 矩阵 R 有如下的分块表示:

$$R = \begin{bmatrix} R_{11} & R_{12} \\ O & R_{22} \end{bmatrix} \begin{matrix} \ell \\ m - \ell \end{matrix}.$$

$\ell \quad m - \ell$

事实上, 若 R 的第 ℓ 行正好对应于它的一个 2×2 阶对角块的第 1 行, 则上述分块就不存在, 此时可取 $\ell - 1$ 作为 ℓ .

再对 V_m 作相应的分块, 即

$$V_m = \begin{bmatrix} V_{m1} & V_{m2} \end{bmatrix}.$$

$\ell \quad m - \ell$

由于 z 有式 (7.138) 的形状, 故式 (7.135) 中的正交矩阵就可取作如下形式:

$$Q = \text{diag}(I_\ell, Q_1),$$

式中: Q_1 为 $m - \ell$ 阶正交矩阵. 对于这样的 Q , 有

$$\widehat{V}_m = V_m Q = \begin{bmatrix} V_{m1} & V_{m2} Q_1 \end{bmatrix},$$

$$\widehat{H}_m = Q^T R Q = \begin{bmatrix} R_{11} & R_{12} Q_1 \\ O & Q_1^T R_{22} Q_1 \end{bmatrix}.$$

这表明, 在进行这一步变换时, V_m 的前 ℓ 列和 R 的 ℓ 阶顺序主子阵并没有改变. 在 Arnoldi 分解式 (7.136) 中 \widehat{V}_m 的前 ℓ 列 V_{m1} 就是 A 的一个不变子空间的一组标准正交基, \widehat{H}_m 的 ℓ 阶顺序主子阵 R_{11} 的特征值就全部是 A 的特征值, 而且后继部分的计算再也不会涉及 V_{m1} 和 R_{11} , 因此就称它们已经被锁定. 随着迭代的进行, 锁定的部分会越来越大, 最终达成收敛.

综合上面的讨论, 可得如下的隐式重开始 Arnoldi 算法.

算法 7.18 (隐式重开始 Arnoldi 方法) 给定 n 阶实矩阵 A , 正整数 k 和 \widehat{m} ($\widehat{m} < k$). 本算法计算 A 的 \widehat{m} 个模最大的特征值和对应的特征向量.

步 1, 初始化

(1) 选择初始向量 v , 并令 $v_1 = v / \|v\|_2$; $m = \widehat{m} + 2$.

(2) 应用算法 2.12 于数据 (A, v_1, k) 上, 产生一个长度为 k 的 Arnoldi 分解

$$AV_k = V_k H_k + \beta_k v_{k+1} e_k^T.$$

步 2, 收敛性判定

(1) 计算 H_k 的实 Schur 分解: $H_k = URU^T$, 其中 U 是正交矩阵, R 是实 Schur 标准形.

(2) 若 H_k^T 的 m 个模最大的特征值已经收敛, 则计算相应 Ritz 向量, 停算; 否则, 进行下一步.

步 3, 重新开始 Arnoldi 过程

(1) 利用 Krylov-Schur 重新开始方法产生一个长度为 m 的 Arnoldi 分解

$$AV_m = V_m H_m + \beta_m v_{m+1} e_m^T.$$

(2) 利用 Arnoldi 过程将该分解扩展为一个长度为 k 的 Arnoldi 分解

$$AV_k = V_k H_k + \beta_k v_{k+1} e_k^T.$$

然后转步 2.

注 7.15 注意到算法 7.18 只能计算 A 的几个模最大的特征值和对应的特征向量. 若要计算 A 在某个数 μ 附近的几个特征值和对应的特征向量, 通常是采取位移求逆的方法, 即将算法 7.18 应用到矩阵 $(A - \mu I)^{-1}$ 上. 一旦 $(A - \mu I)^{-1}$ 的几个模最大的特征值 $\hat{\lambda}_1, \dots, \hat{\lambda}_m$ 已经求得, 则

$$\lambda_i = \mu + 1/\hat{\lambda}_i, \quad i = 1, 2, \dots, m$$

就是 A 的最靠近 μ 的几个特征值.

需注意的是, 应用算法 7.18 于 $(A - \mu I)^{-1}$ 上时, 需要计算形如

$$y = (A - \mu I)^{-1} x$$

的矩阵乘向量, 即需要求解形如

$$(A - \mu I)y = x \tag{7.139}$$

的线性方程组. 由于 Arnoldi 方法对 y 的误差十分敏感, 因此通常采用稀疏列主元的 LU 分解来求解方程组 (7.139). 当然, 在整个迭代过程中只需要作一次分解就够了.

7.6.4 Jacobi-Davidson 方法

7.6.1 节中的 Rayleigh-Ritz 投影方法其实只是子空间迭代方法的一般框架, 它并没有给出子空间的具体选取方法. 7.6.2 节的 Lanczos 方法 (矩阵 A 对称) 和 7.6.3 节的 Arnoldi 方法 (矩阵 A 非对称) 其实都是将投影子空间取为 Krylov 子空间 $\mathcal{K}_k(A, v)$ 的

Rayleigh-Ritz 投影法. 本节将给出投影子空间的另一个取法, 相应的方法称为 Jacobi-Davidson 方法.

设 $A \in \mathbb{C}^{n \times n}$, $\mathcal{K} = \text{span}\{v_1, v_2, \dots, v_m\} \subset \mathbb{C}^n$, $\{v_i\}_{i=1}^m$ 为 \mathcal{K} 中的标准正交基. 记 $V_m = [v_1, v_2, \dots, v_m]$. 选择 $\mu \in \mathbb{C}$ 和 $x \in \mathcal{K}$ 使得

$$(Ax - \mu x) \perp \mathcal{K}.$$

上述条件可以等价地表述为

$$(Ax - \mu x, v) = 0, \quad \forall v \in \mathcal{K}.$$

注意到 x 可表示为 $x = V_m y$, 其中 $y \in \mathbb{C}^m$, 故上述条件又可以等价地表示为

$$V_m^H (AV_m y - \mu V_m y) = 0,$$

即

$$A_m y = \mu y,$$

其中 $A_m = V_m^H A V_m$ 正好是 A 关于 V_m 的 Rayleigh 商. 设 (μ, y) ($\|y\|_2 = 1$) 是 A_m 的特征对, $(\mu, V_m y)$ 是 A 的 Ritz 对, 它是 A 的某个特征对的近似. 若记 $u = V_m y$, 则显然有 $\|u\|_2 = 1$. 现取向量 $z \in \mathbb{C}^n$ 满足 $z^H u = 0$, 且

$$A(u + z) = \lambda(u + z) \iff (A - \lambda I)z = -(A - \lambda I)u. \quad (7.140)$$

上式表明 A 的近似特征向量 u 加上一个和它正交的向量 z 后变为 A 的特征向量. 这启发把上述方程限制在与 u 正交的子空间上, 得

$$(I - uu^H)(A - \lambda I)(I - uu^H)z = -(A - \mu I)u. \quad (7.141)$$

事实上, 考虑到

$$u^H A u = (V_m y)^H A (V_m y) = y^H A_m y = \mu, \quad (7.142)$$

u 与 z 正交以及式 (7.140), 有

$$\begin{aligned} & (I - uu^H)(A - \lambda I)(I - uu^H)z \\ &= (I - uu^H)(A - \lambda I)z = -(I - uu^H)(A - \lambda I)u \\ &= -(A - \lambda I)u + (uu^H A u - \lambda u) = -(A - \mu I)u. \end{aligned}$$

由于式 (7.141) 中的 λ 未知, 用 μ 近似它, 得

$$\tilde{A}z = -\tilde{r}, \quad z^H u = 0, \quad (7.143)$$

式中:

$$\tilde{A} = (I - uu^H)(A - \mu I)(I - uu^H)$$

为 A 在和 u 正交的子空间上的限制; $\tilde{r} = (A - \mu I)u$ 表示 Ritz 对 (μ, u) 的残差. 由式 (7.142) 可知

$$u^H \tilde{r} = u^H (A - \mu I)u = 0,$$

即 $\tilde{\mathbf{r}} \in \text{span}\{\mathbf{u}\}^\perp = \mathcal{R}(\mathbf{I} - \mathbf{u}\mathbf{u}^H)$. 式 (7.143) 也可表示为

$$(\mathbf{A} - \mu\mathbf{I})\mathbf{z} = -\tilde{\mathbf{r}} + \alpha\mathbf{u}, \quad \alpha = \mathbf{u}^H(\mathbf{A} - \mu\mathbf{I})\mathbf{z}, \quad \mathbf{z}^H\mathbf{u} = 0.$$

设 μ 不是 \mathbf{A} 的特征值, 则有

$$\mathbf{z} = (\mathbf{A} - \mu\mathbf{I})^{-1}(-\tilde{\mathbf{r}} + \alpha\mathbf{u}). \quad (7.144)$$

根据求出的 \mathbf{z} , 得

$$\mathbf{u}^H(\mathbf{A} - \mu\mathbf{I})\mathbf{z} = \mathbf{u}^H(-\tilde{\mathbf{r}} + \alpha\mathbf{u}) = \alpha.$$

注意到 \mathbf{u} 与 \mathbf{z} 正交, 则由式 (7.144), 有

$$\alpha = \frac{\mathbf{u}^H(\mathbf{A} - \mu\mathbf{I})^{-1}\tilde{\mathbf{r}}}{\mathbf{u}^H(\mathbf{A} - \mu\mathbf{I})^{-1}\mathbf{u}}.$$

这表明, 当 μ 不是 \mathbf{A} 的特征值时, 方程组 (7.143) 有唯一解.

上面讨论了如何从一个子空间 $\text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m\}$ 出发求 Ritz 向量 \mathbf{u} , 并把它进行校正写为 $\mathbf{u} + \mathbf{z}$, 这里 \mathbf{u} 在该子空间中. 现在利用 \mathbf{z} 构造一个“更大”的子空间, 取

$$\mathbf{v}_{m+1} = \mathbf{z} - \sum_{i=1}^m (\mathbf{z}^H \mathbf{v}_i) \mathbf{v}_i$$

的单位化向量. 下面给出 Jacobi-Davidson 方法的详细算法步骤.

算法 7.19 (Jacobi-Davidson 方法) 给定 $\mathbf{A} \in \mathbb{C}^{n \times n}$, $\mathbf{v}_1 \in \mathbb{C}^n$ 满足 $\|\mathbf{v}_1\|_2 = 1$.

For $m = 1, 2, \dots$ **Do**

(1) 根据 $\mathbf{V}_m = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m]$ 产生 $\mathbf{A}_m = \mathbf{V}_m^H \mathbf{A} \mathbf{V}_m \in \mathbb{C}^{m \times m}$.

(2) 求 \mathbf{A}_m 的某个特征对 (μ, \mathbf{y}) . 令 $\mathbf{u} = \mathbf{V}_m \mathbf{y}$.

(3) 计算 $\tilde{\mathbf{r}} = \mathbf{A}\mathbf{u} - \mu\mathbf{u}$. 若 $\|\tilde{\mathbf{r}}\|_2 \leq \varepsilon$, 停算.

(4) 解方程组 (7.143) 得到 \mathbf{z} .

(5) 取 $\tilde{\mathbf{v}}_{m+1} := \mathbf{z} - \sum_{i=1}^m (\mathbf{z}^H \mathbf{v}_i) \mathbf{v}_i$, $\mathbf{v}_{m+1} = \tilde{\mathbf{v}}_{m+1} / \|\tilde{\mathbf{v}}_{m+1}\|_2$.

Enddo

注 7.16 注意到

$$\begin{aligned} \mathbf{A}_{m+1} &= \mathbf{V}_{m+1}^H \mathbf{A} \mathbf{V}_{m+1} = \begin{bmatrix} \mathbf{V}_m^H \\ \mathbf{v}_{m+1}^H \end{bmatrix} \mathbf{A} [\mathbf{V}_m, \mathbf{v}_{m+1}] \\ &= \begin{bmatrix} \mathbf{A}_m & \mathbf{V}_m^H \mathbf{A} \mathbf{v}_{m+1} \\ \mathbf{v}_{m+1}^H \mathbf{A} \mathbf{V}_m & \mathbf{v}_{m+1}^H \mathbf{A} \mathbf{v}_{m+1} \end{bmatrix}, \end{aligned}$$

故算法 7.19 在第 $m+1$ 次循环时, 只需计算 $\mathbf{V}_m^H \mathbf{A} \mathbf{v}_{m+1}$, $\mathbf{v}_{m+1}^H \mathbf{A} \mathbf{V}_m$ 和 $\mathbf{v}_{m+1}^H \mathbf{A} \mathbf{v}_{m+1}$ 便产生了 \mathbf{A}_{m+1} .

注 7.17 算法 7.19 的计算量主要集中在第 4 步解方程组 (7.143). 可以使用子空间迭代法 GMRES 方法或预处理 GMRES 方法来求解.

算法 7.19 的 MATLAB 程序如下:

```
%Jacobi-Davidson方法程序-jacobidavidson.m
function [mu,u,Vm]=jacobidavidson(A,v,tol,max_it)
%用Jacobi-Davidson求矩阵A的模最大的特征值及相应的特征向量
%输入:A为n阶实方阵,v为初始向量,tol为容许误差,max_it为子空间的最大维数
%输出:mu返回A的模最大特征值,u为相应的特征向量,Vm为子空间基矩阵
if nargin<4, max_it=ceil(size(A,1)/2); end
if nargin<3, tol=1e-6; end
n=size(A,1); m=0; Vm=[ ]; I=eye(n);
x0=rand(n,1); v=v/norm(v);
Vm=[Vm,v]; Am=Vm'*A*Vm;
while(m<max_it)
    m=m+1;
    [Y,D]=eig(Am); [s,t]=max(abs(diag(D)));
    mu=D(t,t); u=Vm*Y(:,t); r=A*u-mu*u;
    if (norm(r)<=tol), break; end
    At=(I-u*u')*(A-mu*I)*(I-u*u');
    [z]=gmres(At,-r); %系统自带的GMRES函数
    v=z; %正交化
    for i=1:m
        v=v-(z'*Vm(:,i))*Vm(:,i);
    end
    v=v/norm(v);
    Am=[Am,Vm'*A*v;v'*A*Vm,v'*A*v];
    Vm=[Vm,v];
end
```

例 7.15 用算法 7.19 计算矩阵

$$A = \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & 2 & 2 & \cdots & 2 \\ 1 & 2 & 3 & \cdots & 3 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 2 & 3 & \cdots & n \end{bmatrix}$$

的模最大特征值. 取 $n = 1000$.

解 编写 MATLAB 脚本程序 ex715.m, 然后命令窗口运行该程序可得到相应的结果. 当投影子空间的维数增加到 $m = 4$ 时, 迭代已满足终止条件, 得到近似特征值 $\tilde{\lambda}$ 和

对应的特征向量 \mathbf{u} 满足 $\|\mathbf{A}\mathbf{u} - \tilde{\lambda}\mathbf{u}\|_2 = 5.3620 \times 10^{-9}$.

假设希望计算矩阵 \mathbf{A} 的某个与 α 最接近的特征值, 可以考虑求 $(\mathbf{A} - \alpha\mathbf{I})^{-1}$ 模最大的特征值 μ , 那么 $\tilde{\lambda} = \frac{1}{\mu} + \alpha$ 即为 \mathbf{A} 的某个与 α 最接近的特征值. 也就是说, 只需用 $\mathbf{B} = (\mathbf{A} - \alpha\mathbf{I})^{-1}$ 去调用算法 7.19 的程序 jacobidavidson.m 即可. 例如, 要计算例 7.15 中的矩阵与 $\alpha = 18$ 最接近的特征值和相应的特征向量, 可以编写如下 MATLAB 脚本程序:

```
tic; n=1000; A=zeros(n,n);
for i=1:n,
    for j=1:i
        A(i,j)=j; A(j,i)=A(i,j);
    end
end
v=rand(n,1); I=eye(n);
alpha=18; B=(A-alpha*I)\I;
[mu,u,Vm]=jacobidavidson(B,v);
lambda=alpha+1/mu,
err1=norm(A*u-lambda*u),d=eig(A);
err2=norm(lambda-d(925)),toc
```

执行上述程序语句, 得到 $\lambda = 17.8762$, $\|\mathbf{A}\mathbf{u} - \lambda\mathbf{u}\|_2 = 8.1153 \times 10^{-9}$.

习题 7

7.1 试证明: 定理 7.6 中的集合 $\bigcup_{i \neq j} \Omega_{ij}$ 是定理 7.5 中的集合 $\bigcup_{i=1}^n G_i$ 的子集.

7.3 设 $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{n \times m}$. 试证明: 矩阵 \mathbf{AB} 与 \mathbf{BA} 具有相同的非零特征值.

7.3 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\xi \pm i\eta$ 为 \mathbf{A} 的一对共轭特征值, 对应的特征向量为 $\mathbf{x} \pm i\mathbf{y}$, 其中 $\xi, \eta \in \mathbb{R}$, $\eta \neq 0$, $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$. 试证明:

- (1) \mathbf{x} 与 \mathbf{y} 线性无关;
- (2) 存在正交矩阵 $\mathbf{Q} \in \mathbb{R}^{n \times n}$, 使满足

$$\mathbf{Q}^T \mathbf{A} \mathbf{Q} = \begin{bmatrix} \mathbf{P} & \mathbf{S} \\ \mathbf{O} & \mathbf{T} \end{bmatrix},$$

其中 $\mathbf{P} \in \mathbb{R}^{2 \times 2}$.

7.4 矩阵的任一特征值及其相应的特征向量称为矩阵的一个特征对. 设 (λ, \mathbf{v}) 是矩阵 \mathbf{A} 的特征对, 试证明:

- (1) 对于任意的常数 α , 则 $(\lambda - \alpha, \mathbf{v})$ 是矩阵 $\mathbf{A} - \alpha\mathbf{I}$ 的特征对;

- (2) 若 $\lambda \neq 0$, 则 $(1/\lambda, \mathbf{v})$ 是矩阵 \mathbf{A}^{-1} 的特征对;
 (3) 若 $\alpha \neq \lambda$, 则 $(1/(\lambda - \alpha), \mathbf{v})$ 是矩阵 $(\mathbf{A} - \alpha \mathbf{I})^{-1}$ 的特征对.

7.5 设 $\mathbf{A} \in \mathbb{C}^{n \times n}$, 对于给定的非零向量 $\mathbf{x} \in \mathbb{C}^n$, 定义

$$\rho(\mathbf{x}) = \frac{\mathbf{x}^H \mathbf{A} \mathbf{x}}{\mathbf{x}^H \mathbf{x}},$$

称之为 \mathbf{x} 对 \mathbf{A} 的 Rayleigh 商. 试证明: 对任意的非零向量 $\mathbf{x} \in \mathbb{C}^n$, 有

$$\|\mathbf{A} \mathbf{x} - \rho(\mathbf{x}) \mathbf{x}\|_2 = \inf_{\mu \in \mathbb{C}} \|\mathbf{A} \mathbf{x} - \mu \mathbf{x}\|_2.$$

7.6 设 \mathbf{A} 为奇异的不可约上 Hessenberg 矩阵. 试证明: 进行一次基本的 QR 迭代后, \mathbf{A} 的零特征值将出现.

7.7 设 $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$, $\mathbf{C} = \mathbf{A} + i\mathbf{B}$, $i = \sqrt{-1}$, $\mathbf{M} = \begin{bmatrix} \mathbf{A} & -\mathbf{B} \\ \mathbf{B} & \mathbf{A} \end{bmatrix}$.

- (1) 证明 \mathbf{C} 为 Hermite 矩阵的充分必要条件是 \mathbf{M} 为实对称矩阵;
 (2) 指出 \mathbf{C} 的特征值和特征向量与 \mathbf{M} 的特征值和特征向量的关系.

7.8 设 \mathbf{A} 为实对称矩阵, $\{\mathbf{A}_k\}$ 是按算法 7.3 (Jacobi 方法) 产生的矩阵序列, 记

$$S(\mathbf{A}_k) = \sum_{\substack{i,j=1 \\ i \neq j}}^n [a_{ij}^{(k)}]^2,$$

证明:

$$\lim_{k \rightarrow \infty} S(\mathbf{A}_k) = 0.$$

7.9 设

$$\mathbf{A} = \begin{bmatrix} (\mathbf{A}_{11})_{3 \times 3} & \mathbf{O}_{3 \times 2} \\ \mathbf{O}_{2 \times 3} & (\mathbf{A}_{22})_{2 \times 2} \end{bmatrix}.$$

又设 λ_i 为 \mathbf{A}_{11} 的特征值, λ_j 为 \mathbf{A}_{22} 的特征值, $\mathbf{x}_i = (\alpha_1, \alpha_2, \alpha_3)^T$ 是 \mathbf{A}_{11} 对应于 λ_i 的特征向量, $\mathbf{y}_j = (\beta_1, \beta_2)^T$ 是 \mathbf{A}_{22} 对应于 λ_j 的特征向量. 求证:

- (1) λ_i 和 λ_j 是 \mathbf{A} 的特征值;
 (2) $\bar{\mathbf{x}}_i = (\alpha_1, \alpha_2, \alpha_3, 0, 0)^T$ 是 \mathbf{A} 对应于 λ_i 的特征向量, $\bar{\mathbf{y}}_j = (0, 0, 0, \beta_1, \beta_2)^T$ 是 \mathbf{A} 对应于 λ_j 的特征向量.

7.10 设 \mathbf{A} 为不可约对称三对角矩阵, 对其实行 QR 方法时, 每个迭代矩阵 $\mathbf{A}_k = \mathbf{Q}_k \mathbf{R}_k$, $k = 0, 1, \dots$, 其中 $\mathbf{A}_0 = \mathbf{A}$. 试证明:

- (1) $\mathbf{R}_k = (r_{ij})$ 为上三角矩阵, 且只有 r_{ij} , $i \leq j \leq i+2$, $i, j = 1, 2, \dots, n$ 可能不为零;
 (2) \mathbf{Q}_k 为不可约上 Hessenberg 矩阵;
 (3) \mathbf{A}_{k+1} 仍为不可约对称三对角矩阵.

7.11 设三对角矩阵

$$A = \begin{bmatrix} a_1 & c_1 & & & \\ b_1 & a_2 & c_2 & & \\ & \ddots & \ddots & \ddots & \\ & & b_{n-2} & a_{n-1} & c_{n-1} \\ & & & b_{n-1} & a_n \end{bmatrix} \in \mathbb{R}^{n \times n},$$

其中 $b_i c_i > 0, i = 1, 2, \dots, n-1$. 试证明: 存在非奇异的对角矩阵 D , 使得 $D^{-1}AD$ 为实对称三对角矩阵.

7.12 设 $A \in \mathbb{C}^{n \times n}$ 有实特征值且满足 $\lambda_1 > \lambda_2 \geq \dots \geq \lambda_{n-1} > \lambda_n$. 现应用幂法于矩阵 $A - \mu I$, 试证明: 选择 $\mu = \frac{1}{2}(\lambda_2 + \lambda_n)$ 时, 所产生的向量序列收敛到属于 λ_1 的特征向量的速度最快.

7.13 设 $A \in \mathbb{C}^{n \times n}, x \in \mathbb{C}^n$, 且 $X = [x, Ax, \dots, A^{n-1}x]$. 试证明: 若 X 非奇异, 则 $X^{-1}AX$ 为上 Hessenberg 矩阵.

7.14 设 A 是上 Hessenberg 矩阵, 其列主元 LU 分解为 $PA = LU$, 这里 P 是置换矩阵, L 为单位下三角矩阵, U 为上三角矩阵. 试证明: $\tilde{A} = U(P^T L)$ 仍为上 Hessenberg 矩阵, 并且相似于 A .

7.15 若给定 $A_0 := A$, 并由

$$A_k - \mu_k I = Q_k R_k, \quad A_{k+1} = R_k Q_k + \mu_k I$$

产生矩阵序列 $\{A_k\}$. 试证明:

$$(Q_0 Q_1 \cdots Q_l)(R_l \cdots R_1 R_0) = (A - \mu_0 I)(A - \mu_1 I) \cdots (A - \mu_l I).$$

参考文献

- [1] 孙继广. 矩阵扰动分析. 北京: 科学出版社, 1987.
- [2] 蔡大用. 数值代数. 北京: 清华大学出版社, 1987.
- [3] 胡家骥. 线性代数方程组的迭代解法. 北京: 科学出版社, 1991.
- [4] 何旭初, 孙文瑜. 广义逆矩阵引论. 南京: 江苏科学技术出版社, 1991.
- [5] 王国荣. 矩阵及算子广义逆. 北京: 科学出版社, 1994.
- [6] 徐树方. 矩阵计算的理论与方法. 北京: 北京大学出版社, 1995.
- [7] 曹志浩. 数值线性代数. 上海: 复旦大学出版社, 1996.
- [8] 戴华. 矩阵论. 北京: 科学出版社, 2001.
- [9] 曹志浩. 变分迭代法. 北京: 科学出版社, 2005.
- [10] 张凯院, 徐仲. 数值代数. 2版. 北京: 科学出版社, 2006.
- [11] 魏木生. 广义最小二乘问题的理论与计算. 北京: 科学出版社, 2006.
- [12] 蒋尔雄. 矩阵计算. 北京: 科学出版社, 2008.
- [13] 李大明. 数值线性代数. 北京: 清华大学出版社, 2010.
- [14] 徐树方. 控制论中的矩阵计算. 北京: 高等教育出版社, 2011.
- [15] 徐树方, 高立, 张平文. 数值线性代数. 2版. 北京: 北京大学出版社, 2013.
- [16] 徐仲, 陆全, 张凯院, 等. H-类矩阵的理论与应用. 北京: 科学出版社, 2013.
- [17] 徐树方, 钱江. 矩阵计算六讲. 北京: 高等教育出版社, 2014.
- [18] 谷同祥, 安恒斌, 刘兴平, 等. 迭代方法和预处理技术(上册). 北京: 科学出版社, 2015.
- [19] 谷同祥, 徐小文, 刘兴平, 等. 迭代方法和预处理技术(下册). 北京: 科学出版社, 2016.
- [20] Golub G H, Van Loan C F. Matrix Computations. 3 ed (中译本). 袁亚湘, 等译. 北京: 人民邮电出版社, 2011.
- [21] Bai Z Z, Golub G H, Ng M K. Hermitian and skew-Hermitian splitting methods for non-Hermitian positive definite linear systems. SIAM J. Matrix Anal. Appl., 2003, 24(3): 603–626.
- [22] Bai Z Z, Golub G H, Pan J Y. Preconditioned Hermitian and skew-Hermitian splitting methods for non-Hermitian positive semidefinite linear systems. Numer. Math., 2004, 98(1): 1–32.
- [23] Niethammer W, Pillis J de, Varga R S. Convergence of block iterative methods applied to sparse least-squares problems. Linear Algebra Appl., 1984, 58(1): 327–341.
- [24] Cline A K, Moler C B, Stewart G W, Wilkinson J H. An estimate for the condition number of a matrix. SIAM J. Numer. Anal., 1979, 16(2): 368–375.
- [25] Freund R W, Golub G H, Nachtigal N M. Iterative solution of linear systems. Acta Numer., 1992, 1(1): 57–100.
- [26] Paige C C, Saunders M A. Solution of sparse indefinite systems of linear equations. SIAM J. Numer. Anal., 1975, 12(4): 617–629.
- [27] Varga R S, Gillis J. Matrix Iterative Analysis. 2 ed. Springer, 2000.
- [28] Young D M. Iterative Solution of Large Linear Systems. Dover Publications, 2003.
- [29] Saad Y. Iterative Methods for Sparse Linear Systems. 2 ed. SIAM, 2003.