

数 值 分 析

第 4 版

李庆扬 王能超 易大义 编

清华大学出版社
施普林格出版社

(京)新登字 158 号

内 容 提 要

本书是为理工科大学各专业普遍开设的“数值分析”课程编写的教材。其内容包括插值与逼近,数值微分与数值积分,非线性方程与线性方程组的数值解法,矩阵的特征值与特征向量计算,常微分方程数值解法。每章附有习题并在书末有部分答案,书末还附有计算实习题和并行算法简介。全书阐述严谨,脉络分明,深入浅出,便于教学。

本书也可作为理工科大学各专业研究生学位课程的教材,并可供从事科学计算的科技工作者参考。

书 名: 数值分析(第 4 版)

作 者: 李庆扬 王能超 易大义 编

出版者: 清华大学出版社 施普林格出版社
(北京清华大学学研大厦,邮编 100084)

<http://www.tup.tsinghua.edu.cn>

印刷者: 北京振华印刷厂

发行者: 新华书店总店北京发行所

开 本: 850 × 1168 1/32 印张: 13.125 字数: 328 千字

版 次: 2001 年 8 月第 4 版 2001 年 8 月第 1 次印刷

书 号: ISBN 7-302-04561-5/O · 259

印 数: 0001 ~ 5000

定 价: 16.00 元

第四版前言

本书由华中理工大学出版社出版至今已 20 年,重新修订的第三版也已 15 年了,印数已近 20 万册,1988 年获国家教委优秀教材二等奖,表明本教材在国内是受欢迎的,仍有存在的价值。为使本书适应新世纪的要求,我们认为对本书重新进行修改是完全必要的。这次修改除保留本书原有风格和基本内容外,修改的原则和内容有以下几点:

(1) 随着计算机技术的发展和普及,数值分析的原理与方法在各学科中的应用越来越多。因此,我们将原来主要面向应用数学专业扩大为面向理工科大学中对数学要求较高的专业的本科生,同时也兼顾到一些院校为各专业研究生开设的“数值分析”学位课程。

(2) 由于科学及计算机的发展,计算机算法语言的多样化及数学软件的普及,要求“数值分析”课程更强调算法原理及理论分析,而对具体算法及编程已有现成数学软件,如 Matlab 等,方便了读者使用。因此,我们对某些算法做了精简,另外也删去了一些较少使用的算法,增加一些实际应用中较重要的内容,如帕德逼近,解线性方程组的 QR 方法及超定方程组最小二乘解,非线性方程组求解的牛顿法,解刚性常微分方程的基本概念等。考虑到很多高校配备了大型多处理机,具备了进行并行计算的条件,故增加了“并行算法及其基本概念”的附录,便于需要进行并行计算的读者对此有初步的了解。

(3) 学习本课程仍应加强上机计算实习,为此,新版增加了计算实习的题目,便于教学,教师可根据实际条件让学生选做其中的

3~5题.由于计算机算法语言发展很快,故不规定用哪种算法语言,目前我们向读者推荐的是集成化软件包 Matlab.

(4) 统一协调,改正错误.本书第三版存在一些不协调之处和各种错误.为保证新版质量,由李庆扬负责对全书整理加工,统一规格并改正旧版中的各种错误.

作者将新版“数值分析”交清华大学出版社重新出版,出版社委派曾多次使用本书的计算数学博士刘颖负责编辑加工,他不但改正了本书的一些错误并对本书修改提出了宝贵意见,提高了本书新版的质量,出版社还在较短时间使本书新版在开学前与读者见面,我们对清华大学出版社及刘颖博士表示衷心感谢.

作 者

2001年5月

第三版说明

本书自 1981 年问世以来, 为许多工科院校所采用, 已先后出过两版, 总发行量达四万余册。1985 年 5 月召开的工科院校计算数学教材评议会(南北会议)确认本书“基本符合应用数学专业的要求, 可作为数值分析课程的教材, 建议作者加以修改后重新出版”。我们遵照这次会议的建议和要求再次进行了修订。新书在出版质量上有了显著的提高。编者诚挚地感谢华中工学院出版社的同志们, 为本书的重版付出了辛勤的劳动。

编 者

1986 年 12 月

第二版前言

1980年7月在大连召开的工科院校“应用数学专业教学学术会议”，根据教育部直属工科院校“应用数学专业教学计划”制定了“数值分析”课大纲，并决定由清华大学、华中工学院、浙江大学合编试用教材。本书就是根据这次会议的决定编写的。全书共分九章，第一、二、三章由李庆扬编写，第四、五、六章由王能超编写，第七、八、九章由易大义编写。

1981年元月在杭州召开的工科院校计算数学第一次教材审稿会，对本教材初稿进行了审查，1982年元月在上海交大召开的第二次计算数学教材审稿会，又对本书第一版提出了修改意见。会议考虑到理工科院校各专业普遍开设“数值分析”课的情况，重新修订了大纲(72学时)。本书第二版就是根据新大纲的要求修改的，它保持了第一版的主要内容及特点，但选材更注意基本要求，减少了部分内容，增加了部分习题答案。本书可作为理工科院校应用数学、力学、物理、计算机软件等专业大学生及其他专业研究生“数值分析”(或“计算方法”)课的教材，也可供学习“计算方法”的科技工作者参考。

我们对参加两次审稿会的同志表示衷心感谢，他们以认真负责的态度对本书提出了许多宝贵意见，对提高教材质量起了很大作用。

编 者
1982年7月

目 录

第 1 章 绪论	(1)
1.1 数值分析研究对象与特点	(1)
1.2 数值计算的误差	(3)
1.2.1 误差来源与分类(3)	
1.2.2 误差与有效数字(4)	
1.2.3 数值运算的误差估计(8)	
1.3 误差定性分析与避免误差危害	(10)
1.3.1 病态问题与条件数(11)	
1.3.2 算法的数值稳定性(12)	
1.3.3 避免误差危害的若干原则(14)	
评注	(18)
习题	(18)
第 2 章 插值法	(21)
2.1 引言	(21)
2.2 拉格朗日插值	(23)
2.2.1 线性插值与抛物插值(23)	
2.2.2 拉格朗日插值多项式(26)	
2.2.3 插值余项与误差估计(28)	
2.3 均差与牛顿插值公式	(31)
2.3.1 均差及其性质(31)	
2.3.2 牛顿插值公式(33)	
2.4 差分与等距节点插值	(35)
2.4.1 差分及其性质(35)	

目 录

2.4.2 等距节点插值公式(38)	
2.5 埃尔米特插值.....	(41)
2.6 分段低次插值.....	(45)
2.6.1 高次插值的病态性质(45)	
2.6.2 分段线性插值(47)	
2.6.3 分段三次埃尔米特插值(48)	
2.7 三次样条插值.....	(51)
2.7.1 三次样条函数(51)	
2.7.2 样条插值函数的建立(52)	
2.7.3 误差界与收敛性(57)	
评注	(58)
习题	(58)

第3章 函数逼近与曲线拟合	(61)
3.1 函数逼近的基本概念.....	(61)
3.1.1 函数逼近与函数空间(61)	
3.1.2 范数与赋范线性空间(64)	
3.1.3 内积与内积空间(65)	
3.2 正交多项式.....	(69)
3.2.1 正交函数族与正交多项式(69)	
3.2.2 勒让德多项式(71)	
3.2.3 切比雪夫多项式(74)	
3.2.4 其他常用的正交多项式(77)	
3.3 最佳一致逼近多项式.....	(78)
3.3.1 基本概念及其理论(78)	
3.3.2 最佳一次逼近多项式(81)	
3.4 最佳平方逼近.....	(83)
3.4.1 最佳平方逼近及其计算(83)	
3.4.2 用正交函数族作最佳平方逼近(87)	
3.5 曲线拟合的最小二乘法.....	(90)

目 录

3.5.1	最小二乘法及其计算(90)
3.5.2	用正交多项式做最小二乘拟合(96)
3.6	最佳平方三角逼近与快速傅里叶变换..... (99)
3.6.1	最佳平方三角逼近与三角插值(99)
3.6.2	快速傅氏变换(FFT)(102)
3.7	有理逼近 (108)
3.7.1	有理逼近与连分式(108)
3.7.2	帕德逼近(110)
评注 (114)
习题 (115)
 第 4 章 数值积分与数值微分..... (118)	
4.1	引言 (118)
4.1.1	数值求积的基本思想(118)
4.1.2	代数精度的概念(120)
4.1.3	插值型的求积公式(121)
4.1.4	求积公式的收敛性与稳定性(122)
4.2	牛顿-柯特斯公式 (123)
4.2.1	柯特斯系数(123)
4.2.2	偶阶求积公式的代数精度(125)
4.2.3	几种低阶求积公式的余项(126)
4.3	复化求积公式 (127)
4.3.1	复化梯形公式(128)
4.3.2	复化辛普森求积公式(129)
4.4	龙贝格求积公式 (131)
4.4.1	梯形法的递推化(131)
4.4.2	龙贝格算法(133)
4.4.3	理查森外推加速法(135)
4.5	高斯求积公式 (139)
4.5.1	一般理论(139)

目 录

4.5.2 高斯-勒让德求积公式(144)	
4.5.3 高斯-切比雪夫求积公式(146)	
4.6 数值微分	(148)
4.6.1 中点方法与误差分析(148)	
4.6.2 插值型的求导公式(150)	
4.6.3 利用数值积分求导(153)	
4.6.4 三次样条求导(155)	
4.6.5 数值微分的外推算法(156)	
评注	(157)
习题	(158)
 第 5 章 解线性方程组的直接方法	(161)
5.1 引言与预备知识	(161)
5.1.1 引言(161)	
5.1.2 向量和矩阵(162)	
5.1.3 特殊矩阵(163)	
5.2 高斯消去法	(165)
5.2.1 高斯消去法(166)	
5.2.2 矩阵的三角分解(172)	
5.3 高斯主元素消去法	(174)
5.3.1 列主元素消去法(176)	
5.3.2 高斯-若当消去法(180)	
5.4 矩阵三角分解法	(183)
5.4.1 直接三角分解法(183)	
5.4.2 平方根法(188)	
5.4.3 追赶法(193)	
5.5 向量和矩阵的范数	(196)
5.6 误差分析	(205)
5.6.1 矩阵的条件数(205)	
5.6.2 迭代改善法(212)	

目 录

5.7 矩阵的正交三角化及应用	(214)
5.7.1 初等反射阵(215)	
5.7.2 平面旋转矩阵(218)	
5.7.3 矩阵的 QR 分解(220)	
5.7.4 求解超定方程组(225)	

评注.....	(228)
---------	-------

习题.....	(229)
---------	-------

第 6 章 解线性方程组的迭代法.....	(233)
------------------------------	-------

6.1 引言	(233)
--------------	-------

6.2 基本迭代法	(236)
-----------------	-------

6.2.1 雅可比迭代法(237)	
-------------------	--

6.2.2 高斯-塞德尔迭代法(238)	
----------------------	--

6.2.3 解大型稀疏线性方程组的逐次超松弛迭代法(240)	
--------------------------------	--

6.3 迭代法的收敛性	(243)
-------------------	-------

6.3.1 一阶定常迭代法的基本定理(243)	
-------------------------	--

6.3.2 关于解某些特殊方程组迭代法的收敛性(249)	
------------------------------	--

6.4 分块迭代法	(256)
-----------------	-------

评注.....	(259)
---------	-------

习题.....	(259)
---------	-------

第 7 章 非线性方程求根.....	(261)
---------------------------	-------

7.1 方程求根与二分法	(261)
--------------------	-------

7.1.1 引言(261)	
---------------	--

7.1.2 二分法(262)	
----------------	--

7.2 迭代法及其收敛性	(265)
--------------------	-------

7.2.1 不动点迭代法(265)	
-------------------	--

7.2.2 不动点的存在性与迭代法的收敛性(267)	
----------------------------	--

7.2.3 局部收敛性与收敛阶(269)	
----------------------	--

目 录

7.3	迭代收敛的加速方法	(272)
7.3.1	埃特金加速收敛方法(272)	
7.3.2	斯蒂芬森迭代法(273)	
7.4	牛顿法	(276)
7.4.1	牛顿法及其收敛性(276)	
7.4.2	牛顿法应用举例(278)	
7.4.3	简化牛顿法与牛顿下山法(279)	
7.4.4	重根情形(282)	
7.5	弦截法与抛物线法	(283)
7.5.1	弦截法(283)	
7.5.2	抛物线法(285)	
7.6	解非线性方程组的牛顿迭代法	(287)
评注	(289)
习题	(290)

第8章	矩阵特征值问题计算	(292)
8.1	引言	(292)
8.2	幂法及反幂法	(299)
8.2.1	幂法(299)	
8.2.2	加速方法(304)	
8.2.3	反幂法(308)	
8.3	豪斯霍尔德方法	(312)
8.3.1	引言(312)	
8.3.2	用正交相似变换约化一般矩阵为上三角阵(313)	
8.3.3	用正交相似变换约化对称阵为对称三对角阵(317)	
8.4	QR方法	(319)
8.4.1	QR算法(319)	
8.4.2	带原点位移的 QR 方法(322)	
8.4.3	用单步 QR 方法计算上三角阵特征值(325)	

目 录

8.4.4 [*]	双步 QR 方法(隐式 QR 方法)	(329)
评注		(333)
习题		(333)
第 9 章 常微分方程初值问题数值解法		(336)
9.1	引言	(336)
9.2	简单的数值方法与基本概念	(337)
9.2.1	欧拉法与后退欧拉法	(337)
9.2.2	梯形方法	(340)
9.2.3	单步法的局部截断误差与阶	(341)
9.2.4	改进的欧拉公式	(343)
9.3	龙格-库塔方法	(344)
9.3.1	显式龙格-库塔法的一般形式	(344)
9.3.2	二阶显式 R-K 方法	(346)
9.3.3	三阶与四阶显式 R-K 方法	(348)
9.3.4	变步长的龙格-库塔方法	(351)
9.4	单步法的收敛性与稳定性	(352)
9.4.1	收敛性与相容性	(352)
9.4.2	绝对稳定性与绝对稳定域	(355)
9.5	线性多步法	(360)
9.5.1	线性多步法的一般公式	(360)
9.5.2	阿当姆斯显式与隐式公式	(362)
9.5.3	米尔尼方法与辛普森方法	(366)
9.5.4	汉明方法	(367)
9.5.5	预测-校正方法	(368)
9.5.6	构造多步法公式的注记和例	(371)
9.6	方程组和高阶方程	(373)
9.6.1	一阶方程组	(373)
9.6.2	化高阶方程为一阶方程组	(376)
9.6.3	刚性方程组	(378)

目 录

评注	(380)
习题	(381)
计算实习题	(383)
附录 并行算法及其基本概念	(388)
参考文献	(398)
部分习题答案	(400)

第1章 絮 论

1.1 数值分析研究对象与特点

数值分析是计算数学的一个主要部分,计算数学是数学科学的一个分支,它研究用计算机求解各种数学问题的数值计算方法及其理论与软件实现.为了具体说明数值分析的研究对象,我们考察用计算机解决科学计算问题时经历的几个过程:



由实际问题的提出到上机求得问题解答的整个过程都可看作是应用数学的范畴.如果细分的话,由实际问题应用有关科学知识和数学理论建立数学模型这一过程,通常作为应用数学的任务.而根据数学模型提出求解的数值计算方法直到编出程序上机算出结果,这一过程则是计算数学的任务,也是数值分析研究的对象.数值分析的内容包括函数的数值逼近、数值微分与数值积分、非线性方程数值解、数值线性代数、常微和偏微数值解等,它们都是以数学问题为研究对象的,只是它不像纯数学那样只研究数学本身的理论,而是把理论与计算紧密结合,着重研究数学问题的数值方法及其理论.

数值分析也称计算方法,但不应片面地理解为各种数值方法的简单罗列和堆积,同数学分析一样,它也是一门内容丰富,研究方法深刻,有自身理论体系的课程,既有纯数学高度抽象性与严密科学性的特点,又有应用的广泛性与实际试验的高度技术性的特

点,是一门与计算机使用密切结合的实用性很强的数学课程。为了说明它与纯数学课的不同,例如考虑线性方程组数值解,在“线性代数”课程中只介绍解的存在唯一性及有关理论和精确解法,用这些理论和方法还不能在计算机上解上百个未知数的方程组,更不用说求解十几万个未知数的方程组了,要求解这类问题还应根据方程特点,研究适合计算机使用的,满足精度要求,计算时间省的有效算法及其相关的理论。在实现这些算法时往往还要根据计算机的容量、字长、速度等指标,研究具体的求解步骤和程序设计技巧。有的方法在理论上虽不够严格,但通过实际计算、对比分析等手段,证明是行之有效的方法,也应采用。这些就是数值分析具有的特点,概括起来有四点:

第一,面向计算机,要根据计算机特点提供切实可行的有效算法。即算法只能包括加、减、乘、除运算和逻辑运算,这些运算是计算机能直接处理的运算。

第二,有可靠的理论分析,能任意逼近并达到精度要求,对近似算法要保证收敛性和数值稳定性,还要对误差进行分析。这些都建立在相应数学理论的基础上。

第三,要有好的计算复杂性,时间复杂性好是指节省时间,空间复杂性好是指节省存储量,这也是建立算法要研究的问题,它关系到算法能否在计算机上实现。

第四,要有数值实验,即任何一个算法除了从理论上要满足上述三点外,还要通过数值试验证明是行之有效的。

根据“数值分析”课程的特点,学习时我们首先要注意掌握方法的基本原理和思想,要注意方法处理的技巧及其与计算机的结合,要重视误差分析、收敛性及稳定性基本理论;其次,要通过例子,学习使用各种数值方法解决实际计算问题;最后,为了掌握本课的内容,还应做一定数量的理论分析与计算练习。由于本课内容包括了微积分、代数、常微分方程的数值方法,读者必须掌握这

几门课的基本内容才能学好这门课程 .

1 2 数值计算的误差

1 2 1 误差来源与分类

用计算机解决科学计算问题首先要建立数学模型, 它是对被描述的实际问题进行抽象、简化而得到的, 因而是近似的. 我们把数学模型与实际问题之间出现的这种误差称为模型误差. 只有实际问题提法正确, 建立数学模型时又抽象、简化得合理, 才能得到好的结果. 由于这种误差难于用数量表示, 通常都假定数学模型是合理的, 这种误差可忽略不计, 在“数值分析”中不予讨论. 在数学模型中往往还有一些根据观测得到的物理量, 如温度、长度、电压等等, 这些参量显然也包含误差. 这种由观测产生的误差称为观测误差, 在“数值分析”中也不讨论这种误差. 数值分析只研究用数值方法求解数学模型产生的误差.

当数学模型不能得到精确解时, 通常要用数值方法求它的近似解, 其近似解与精确解之间的误差称为截断误差或方法误差. 例如, 函数 $f(x)$ 用泰勒(Taylor)多项式

$$P_n(x) = f(0) + \frac{f'(0)}{1!}x + \frac{f''(0)}{2!}x^2 + \dots + \frac{f^{(n)}(0)}{n!}x^n$$

近似代替, 则数值方法的截断误差是

$$R_n(x) = f(x) - P_n(x) = \frac{f^{(n+1)}(\cdot)}{(n+1)!}x^{n+1}, \quad \text{在 } 0 \text{ 与 } x \text{ 之间}.$$

有了求解数学问题的计算公式以后, 用计算机做数值计算时, 由于计算机的字长有限, 原始数据在计算机上表示会产生误差, 计算过程又可能产生新的误差, 这种误差称为舍入误差. 例如, 用 3.14159 近似代替 π , 产生的误差

$$R = \pi - 3.14159 = 0.0000026\dots$$

就是舍入误差.

此外由原始数据或机器中的十进制数转化为二进制数产生的初始误差对数值计算也将造成影响, 分析初始数据的误差通常也归结为舍入误差.

研究计算结果的误差是否满足精度要求就是误差估计问题, 本书主要讨论算法的截断误差与舍入误差, 而截断误差将结合具体算法讨论. 为分析数值运算的舍入误差, 先要对误差基本概念做简单介绍.

1.2.2 误差与有效数字

定义 1 设 x 为准确值, x^* 为 x 的一个近似值, 称 $e^* = x^* - x$ 为近似值的绝对误差, 简称误差.

注意这样定义的误差 e^* 可正可负, 当绝对误差为正时近似值偏大, 叫强近似值; 当绝对误差为负时近似值偏小, 叫弱近似值.

通常我们不能算出准确值 x , 也不能算出误差 e^* 的准确值, 只能根据测量工具或计算情况估计出误差的绝对值不超过某正数 ϵ^* , 也就是误差绝对值的一个上界. ϵ^* 叫做近似值的误差限, 它总是正数. 例如, 用毫米刻度的米尺测量一长度 x , 读出和该长度接近的刻度 x^* , x^* 是 x 的近似值, 它的误差限是 0.5mm, 于是 $|x^* - x| \leq 0.5\text{mm}$; 如读出的长度为 765mm, 则有 $|765 - x| \leq 0.5$. 从这个不等式我们仍不知道准确的 x 是多少, 但知道 764.5 $\leq x \leq 765.5$, 说明 x 在区间 $[764.5, 765.5]$ 内.

对于一般情形 $|x^* - x| \leq \epsilon^*$, 即

$$x^* - \epsilon^* \leq x \leq x^* + \epsilon^*,$$

这个不等式有时也表示为

$$x = x^* \pm \epsilon^*.$$

误差限的大小还不能完全表示近似值的好坏. 例如, 有两个量 $x = 10 \pm 1$, $y = 1000 \pm 5$, 则

$$x^* = 10, \quad x = 1; \quad y^* = 1000, \quad y = 5.$$

虽然 y^* 比 x^* 大 4 倍, 但 $y^*/y = \frac{5}{1000} = 0.5\%$ 比 $x^*/x = \frac{1}{10} = 10\%$ 要小得多, 这说明 y^* 近似 y 的程度比 x^* 近似 x 的程度要好得多。所以, 除考虑误差的大小外, 还应考虑准确值 x 本身的大

小。我们把近似值的误差 e^* 与准确值 x 的比值

$$\frac{e^*}{x} = \frac{x^* - x}{x}$$

称为近似值 x^* 的相对误差, 记作 e_r^* 。

在实际计算中, 由于真值 x 总是不知道的, 通常取

$$e_r^* = \frac{e^*}{x^*} = \frac{x^* - x}{x^*}$$

作为 x^* 的相对误差, 条件是 $e^* = \frac{e^*}{x^*}$ 较小, 此时

$$\frac{e^*}{x} - \frac{e^*}{x^*} = \frac{e^*(x^* - x)}{x^* x} = \frac{(e^*)^2}{x^*(x^* - e^*)} = \frac{(e^*/x^*)^2}{1 - (e^*/x^*)}$$

是 e_r^* 的平方项级, 故可忽略不计。

相对误差也可正可负, 它的绝对值上界叫做相对误差限, 记作

$$r^*, \text{ 即 } r^* = \frac{|x^*|}{|x|}.$$

根据定义, 上例中 $\frac{|x^*|}{|x|} = 10\%$ 与 $\frac{|y^*|}{|y|} = 0.5\%$ 分别为 x 与 y 的相对误差限, 可见 y^* 近似 y 的程度比 x^* 近似 x 的程度好。

当准确值 x 有多位数时, 常常按四舍五入的原则得到 x 的前几位近似值 x^* , 例如

$$x = \dots = 3.14159265\dots$$

取 3 位 $x_3^* = 3.14, \quad x_3 = 0.002,$

取 5 位 $x_5^* = 3.1416, \quad x_5 = 0.00008,$

它们的误差都不超过末位数字的半个单位, 即

$$/ - 3.14 / \frac{1}{2} \times 10^{-2}, \quad / - 3.1416 / \frac{1}{2} \times 10^{-4}.$$

定义2 若近似值 x^* 的误差限是某一位的半个单位, 该位到 x^* 的第一位非零数字共有 n 位, 就说 x^* 有 n 位有效数字. 它可表示为

$$x^* = \pm 10^m \times (a_0 + a_1 \times 10^{-1} + \dots + a_n \times 10^{-(n-1)}), \quad (2.1)$$

其中 a_i ($i=1, \dots, n$) 是 0 到 9 中的一个数字, $a_0 \neq 0$, m 为整数, 且

$$/ x - x^* / \leq \frac{1}{2} \times 10^{m-n+1}. \quad (2.2)$$

如取 $x^* = 3.14$ 作 x 的近似值, x^* 就有 3 位有效数字, 取 $x^* = 3.1416$, x^* 就有 5 位有效数字.

例1 按四舍五入原则写出下列各数具有 5 位有效数字的近似数: 187.9325, 0.03785551, 8.000033, 2.7182818.

按定义, 上述各数具有 5 位有效数字的近似数分别是

$$187.93, 0.037856, 8.0000, 2.7183.$$

注意 $x = 8.000033$ 的 5 位有效数字近似数是 8.0000 而不是 8, 因为 8 只有 1 位有效数字.

例2 重力常数 g , 如果以 m/s^2 为单位, $g = 9.80 \text{ m/s}^2$; 若以 km/s^2 为单位, $g = 0.00980 \text{ km/s}^2$, 它们都具有 3 位有效数字, 因为按第一种写法

$$/ g - 9.80 / \leq \frac{1}{2} \times 10^{-2},$$

据(2.1), 这里 $m=0$, $n=3$; 按第二种写法

$$/ g - 0.00980 / \leq \frac{1}{2} \times 10^{-5},$$

这里 $m=-3$, $n=3$. 它们虽然写法不同, 但都具有 3 位有效数字.

至于绝对误差限, 由于单位不同结果也不同, $|e|_1^* = \frac{1}{2} \times 10^{-2} \text{ m/s}^2$,

$|e|_2^* = \frac{1}{2} \times 10^{-5} \text{ km/s}^2$, 而相对误差都是

$$r^* = 0.005 / 9.80 = 0.000005 / 0.00980.$$

注意相对误差与绝对误差限是无量纲的, 而绝对误差与误差限是有量纲的.

例 2 说明有效位数与小数点后有多少位数无关. 然而, 从(2.2)可得到具有 n 位有效数字的近似数 x^* , 其绝对误差限为

$$r^* = \frac{1}{2} \times 10^{m-n+1},$$

在 m 相同的情况下, n 越大则 10^{m-n+1} 越小, 故有效位数越多, 绝对误差限越小.

至于有效数字与相对误差限的关系, 有

定理 1 设近似数 x^* 表示为

$$x^* = \pm 10^m \times (a_0 + a_1 \times 10^{-1} + \dots + a_{l-1} \times 10^{-(l-1)}), \quad (2.1)$$

其中 a_i ($i=1, 2, \dots, l$) 是 0 到 9 中的一个数字, $a_0 \neq 0$, m 为整数. 若 x^* 具有 n 位有效数字, 则其相对误差限为

$$r^* = \frac{1}{2a_0} \times 10^{-(n-1)};$$

反之, 若 x^* 的相对误差限 $r^* = \frac{1}{2(a_0 + 1)} \times 10^{-(n-1)}$, 则 x^* 至少具有 n 位有效数字.

证明 由(2.1) 可得

$$a_0 \times 10^m / |x^*| < (a_0 + 1) \times 10^m,$$

当 x^* 有 n 位有效数字时

$$r^* = \frac{|x - x^*|}{|x^*|} = \frac{0.5 \times 10^{m-n+1}}{a_0 \times 10^m} = \frac{1}{2a_0} \times 10^{-n+1};$$

反之, 由

$$\begin{aligned} & |x - x^*| \\ &= |x^*| r^* < (a_0 + 1) \times 10^m \times \frac{1}{2(a_0 + 1)} \times 10^{-n+1} \end{aligned}$$

$$= 0.5 \times 10^{m-n+1},$$

故 x^* 至少有 n 位有效数字. 定理证完.

定理说明, 有效位数越多, 相对误差限越小.

例 3 要使 20 的近似值的相对误差限小于 0.1%, 要取几位有效数字?

设取 n 位有效数字, 由定理 1, $\frac{1}{2a_1} \times 10^{-n+1}$. 由于 $20 =$

4.4..., 知 $a = 4$, 故只要取 $n = 4$, 就有

$$\frac{1}{2 \cdot 4} \times 10^{-3} < 10^{-3} = 0.1\%,$$

即只要对 20 的近似值取 4 位有效数字, 其相对误差限就小于 0.1%. 此时由开方表得 20 = 4.472.

1.2.3 数值运算的误差估计

两个近似数 x_1^* 与 x_2^* , 其误差限分别为 (x_1^*) 及 (x_2^*) , 它们进行加、减、乘、除运算得到的误差限分别为

$$(x_1^* \pm x_2^*) = (x_1^*) + (x_2^*);$$

$$(x_1^* x_2^*) = |x_1^*| (x_2^*) + |x_2^*| (x_1^*);$$

$$(x_1^* / x_2^*) = \frac{|x_1^*| / (x_2^*) + |x_2^*| / (x_1^*)}{|x_2^*|^2} \quad (x_2^* \neq 0).$$

更一般的情况是, 当自变量有误差时计算函数值也产生误差, 其误差限可利用函数的泰勒展开式进行估计. 设 $f(x)$ 是一元函数, x 的近似值为 x^* , 以 $f(x^*)$ 近似 $f(x)$, 其误差界记作 $(f(x^*))$, 可用泰勒展开

$$f(x) - f(x^*) = f(x^*)(x - x^*) + \frac{f''(\cdot)}{2}(x - x^*)^2,$$

介于 x, x^* 之间,

取绝对值得

$$|f(x) - f(x^*)| \leq |f(x^*)| + \frac{|f''(\cdot)|}{2}(x^*)^2.$$

假定 $f(x^*)$ 与 $f(\bar{x})$ 的比值不太大, 可忽略 (x^*) 的高阶项, 于是可得计算函数的误差限

$$(f(x^*)) \quad / \quad f(\bar{x}) \quad / \quad (x^*).$$

当 f 为多元函数时, 例如计算 $A = f(x_1, \dots, x_n)$. 如果 x_1, \dots, x_n 的近似值为 $\bar{x}_1^*, \dots, \bar{x}_n^*$, 则 A 的近似值为 $\bar{A}^* = f(\bar{x}_1^*, \dots, \bar{x}_n^*)$, 于是由泰勒展开得函数值 \bar{A}^* 的误差 $e(\bar{A}^*)$ 为

$$\begin{aligned} e(\bar{A}^*) &= \bar{A}^* - A = f(\bar{x}_1^*, \dots, \bar{x}_n^*) - f(x_1, \dots, x_n) \\ &= \sum_{k=1}^n \frac{f(\bar{x}_1^*, \dots, \bar{x}_{k-1}^*, \bar{x}_k^*, \bar{x}_{k+1}, \dots, \bar{x}_n)}{x_k} (x_k^* - x_k) \\ &= \sum_{k=1}^n \frac{\bar{f}}{x_k} e_k^*, \end{aligned}$$

于是误差限

$$(A^*) \quad \left| \quad \frac{\bar{f}}{x_k} \right|^* \quad |(x_k^*)|; \quad (2.3)$$

而 A^* 的相对误差限为

$$r^* = r(A^*) = \frac{(A^*)}{|A^*|} \quad \left| \quad \frac{\bar{f}}{x_k} \right|^* \quad \left| \frac{(x_k^*)}{|A^*|} \right|. \quad (2.4)$$

例 4 已测得某场地长 l 的值为 $\bar{l}^* = 110\text{m}$, 宽 d 的值为 $\bar{d}^* = 80\text{m}$, 已知 $|l - \bar{l}^*| = 0.2\text{m}$, $|d - \bar{d}^*| = 0.1\text{m}$. 试求面积 $s = ld$ 的绝对误差限与相对误差限.

解 因 $s = ld$, $\frac{s}{l} = d$, $\frac{s}{d} = l$, 由(2.3)知

$$(s^*) \quad \left| \quad \frac{s}{l} \right|^* \quad |(\bar{l}^*)| + \left| \quad \frac{s}{d} \right|^* \quad |(\bar{d}^*)|,$$

其中

$$\frac{s}{l}^* = d^* = 80\text{m}, \quad \frac{s}{d}^* = l^* = 110\text{m},$$

而 $(\bar{l}^*) = 0.2\text{m}$, $(\bar{d}^*) = 0.1\text{m}$,

于是绝对误差限

$$(s^*) \quad 80 \times (0.2) + 110 \times (0.1) = 27(m^2);$$

相对误差限

$$r(s^*) = \frac{(s^*)}{|s^*|} = \frac{(s^*)}{l^* d^*} = \frac{27}{8800} = 0.31\%.$$

1.3 误差定性分析与避免误差危害

数值运算中的误差分析是个很重要而复杂的问题, 上节讨论了不精确数据运算结果的误差限, 它只适用于简单情形, 然而一个工程或科学计算问题往往要运算千万次, 由于每步运算都有误差, 如果每步都做误差分析是不可能的, 也不科学, 因为误差积累有正有负, 绝对值有大有小, 都按最坏情况估计误差限得到的结果比实际误差大得多, 这种保守的误差估计不反映实际误差积累. 考虑到误差分布的随机性, 有人用概率统计方法, 将数据和运算中的舍入误差视为适合某种分布的随机变量, 然后确定计算结果的误差分布, 这样得到的误差估计更接近实际, 这种方法称为概率分析法.

20世纪60年代以后对舍入误差分析提出了一些新方法, 较重要的有以下两种:

1. 向后误差分析法是把新算出的量由某个公式表达, 它仅含基本算术运算, 如假定 a_1, \dots, a_n 是前面已算出的量或原始数据, 新算出量

$$x = g(a_1, \dots, a_n).$$

若 a_i 的摄动为 ε_i , 使得由浮点运算得出结果为

$$x_{fl} = g(a_1 + \varepsilon_1, \dots, a_n + \varepsilon_n).$$

则可根据 ε_i 的界由摄动理论估计最后舍入误差 $|x - x_{fl}|$ 的界, 威克逊(Wilkinson)将这种方法应用于数值代数(矩阵运算)的误差分析, 取得较好效果.

2. 区间分析法是把参加运算的数 x, y, z, \dots 都看成区间量 X, Y, Z, \dots , 根据区间运算规则求得最后结果的近似值及误差限. 例如, x, y 的近似数为 $[x^-, x^+]$, 由于 $|x^- - x^+|$ 是 x 的误差限, $[x^-, x^+] = X$, $y = [y^-, y^+] = Y$, 若计算 $z = x * y$ ($*$ 为运算符号), 由 $Z = X * Y = [z^- z^+] = [z^- - z^+, z^+ + z^-]$, 则 z 为所求近似值, 而 $|z^- - z^+|$ 则为误差限.

上面简略介绍了误差分析的几种方法, 但都不是十分有效的, 目前尚无有效的方法对误差做出定量估计. 为了确保数值计算结果的正确性, 首先需对数值计算问题做定性分析, 为此本节讨论以下三个问题.

1.3.1 病态问题与条件数

对一个数值问题本身如果输入数据有微小扰动(即误差), 引起输出数据(即问题解)相对误差很大, 这就是病态问题. 例如计算函数值 $f(x)$ 时, 若 x 有扰动 $x = x - x^*$, 其相对误差为 $\frac{x}{x}$, 函数值 $f(x^*)$ 的相对误差为 $\frac{f(x) - f(x^*)}{f(x)}$. 相对误差比值

$$\left| \frac{f(x) - f(x^*)}{f(x)} \right| / \left| \frac{x}{x} \right| = \left| \frac{xf(x)}{f(x)} \right| = C_p, \quad (3.1)$$

C_p 称为计算函数值问题的条件数. 自变量相对误差一般不会太大, 如果条件数 C_p 很大, 将引起函数值相对误差很大, 出现这种情况的问题就是病态问题.

例如, $f(x) = x^n$, 则有 $C_p = n$, 它表示相对误差可能放大 n 倍. 如 $n=10$, 有 $f(1) = 1, f(1.02) = 1.24$, 若取 $x=1, x^*=1.02$ 自变量相对误差为 2%, 函数值相对误差为 24%, 这时问题可以认为是病态的. 一般情况条件数 $C_p > 10$ 就认为是病态, C_p 越大病态越严重.

其他计算问题也要分析是否病态. 例如解线性方程组, 如果输入数据有微小误差引起解的巨大误差, 就认为是病态方程组, 我们将在第5章用矩阵的条件数来分析这种现象.

1.3.2 算法的数值稳定性

用一个算法进行计算, 由于初始数据误差在计算中传播使计算结果误差增长很快就是数值不稳定的, 先看下例.

例 5 计算 $I_n = \int_0^1 x^n e^x dx$ ($n = 0, 1, \dots$) 并估计误差.

由分部积分可得计算 I_n 的递推公式

$$I_n = 1 - nI_{n-1} \quad (n = 1, 2, \dots), \quad (3.2)$$

$$I_0 = \int_0^1 e^x dx = 1 - e^{-1}.$$

若计算出 I_0 , 代入(3.2), 可逐次求出 I_1, I_2, \dots 的值. 要算出 I_0 就要先计算 e^{-1} , 若用泰勒多项式展开部分和

$$e^{-1} = 1 + (-1) + \frac{(-1)^2}{2!} + \dots + \frac{(-1)^k}{k!},$$

并取 $k=7$, 用4位小数计算, 则得 $e^{-1} = 0.3679$, 截断误差 $R = |e^{-1} - 0.3679| = \frac{1}{8!} < \frac{1}{4} \times 10^{-4}$. 计算过程中小数点后第5位的数字按四舍五入原则舍入, 由此产生的舍入误差这里先不讨论. 当初值取为 $I_0 = 0.6321$ 时, 用(3.2)递推的计算公式为

$$(A) \quad \begin{aligned} I_0 &= 0.6321; \\ I_n &= 1 - I_{n-1} \quad (n = 1, 2, \dots). \end{aligned}$$

计算结果见表1-1的*璇*列. 用*璇*近似 I_0 产生的误差 $E_0 = I_0 - \text{璇}$ 就是初值误差, 它对后面计算结果是有影响的.

表 1-1

n	ϑ_n (用(A)算)	I^* (用(B)算)	n	ϑ_n (用(A)算)	I^* (用(B)算)
0	0.6321	0.6321	5	0.1480	0.1455
1	0.3679	0.3679	6	0.1120	0.1268
2	0.2642	0.2643	7	0.2160	0.1121
3	0.2074	0.2073	8	-0.7280	0.1035
4	0.1704	0.1708	9	7.552	0.0684

从表中看到 ϑ_n 出现负值, 这与一切 $I_n > 0$ 相矛盾. 实际上, 由积分估值得

$$\begin{aligned} \frac{e^{-1}}{n+1} &= e^{-1} \left(\min_{0 \leq x \leq 1} e^x \right)^{-1} \int_0^n x^n dx < I_n < e^{-1} \left(\max_{0 \leq x \leq 1} e^x \right)^{-1} \int_0^n x^n dx \\ &= \frac{1}{n+1}. \end{aligned} \quad (3.3)$$

因此, 当 n 较大时, 用 ϑ_n 近似 I_n 显然是不正确的. 这里计算公式与每步计算都是正确的, 那么是什么原因使计算结果错误呢? 主要就是初值 ϑ_0 有误差 $E_0 = I_0 - \vartheta_0$, 由此引起以后各步计算的误差 $E_n = I_n - \vartheta_n$ 满足关系

$$E_n = -nE_{n-1} \quad (n = 1, 2, \dots).$$

容易推得

$$E_n = (-1)^n n! E_0,$$

这说明 ϑ_n 有误差 E_0 , 则 ϑ_n 就是 E_0 的 $n!$ 倍误差. 例如, $n=8$, 若 $|E_0| = \frac{1}{2} \times 10^{-4}$, 则 $|E_8| = 8! \times |E_0| > 2$. 这就说明 ϑ_n 完全不能近似 I_8 了. 它表明计算公式(A)是数值不稳定的.

我们现在换一种计算方案. 由(3.3)取 $n=9$, 取

$$\frac{e^{-1}}{10} < I_9 < \frac{1}{10},$$

我们粗略取 $I_9 = \frac{1}{2} \cdot \frac{1}{10} + \frac{e^{-1}}{10} = 0.0684 = I_9$, 然后将公式(3.2)倒

过来算, 即由 I_6^* 算出 $I_8^*, I_7^*, \dots, I_0^*$, 公式为

$$I_9^* = 0.0684,$$

$$(B) \quad I_{n+1}^* = \frac{1}{n}(1 - I_n^*) \quad (n = 9, 8, \dots, 1);$$

计算结果见表 1-1 的 I_n^* 列. 我们发现 I_0^* 与 I_0 的误差不超过 10^{-4} . 记 $E_n^* = I_n - I_n^*$, 则 $|E_0^*| = \frac{1}{n!} |E_n^*|$, E_0^* 比 E_n^* 缩小了 $n!$ 倍, 因此, 尽管 E_0^* 较大, 但由于误差逐步缩小, 故可用 I_n^* 近似 I_n . 反之, 当用方案(A)计算时, 尽管初值 I_0 相当准确, 由于误差传播是逐步扩大的, 因而计算结果不可靠. 此例说明, 数值不稳定的算法是不能使用的.

定义 3 一个算法如果输入数据有误差, 而在计算过程中舍入误差不增长, 则称此算法是数值稳定的, 否则称此算法为不稳定的.

在例 5 中算法(B)是数值稳定的, 而算法(A)是不稳定的. 数值不稳定现象属于误差危害现象, 如何防止误差危害下面将进一步讨论.

1.3.3 避免误差危害的若干原则

数值计算中首先要分清问题是否病态和算法是否数值稳定, 计算时还应尽量避免误差危害, 防止有效数字的损失, 下面给出若干原则.

1. 要避免除数绝对值远远小于被除数绝对值的除法

用绝对值小的数作除数舍入误差会增大, 如计算 $\frac{x}{y}$, 若 $0 < |y| \ll |x|$, 则可能对计算结果带来严重影响, 应尽量避免.

例 6 线性方程组

$$0.00001x_1 + x_2 = 1,$$

$$2x_1 + x_2 = 2.$$

的准确解为

$$x_1 = \frac{200000}{399999} = 0.50000125, \quad x_2 = \frac{199998}{199999} = 0.999995.$$

现在四位浮点十进制数(仿机器实际计算)下用消去法求解, 上述方程写成

$$\begin{aligned} 10^{-4} \cdot 0.1000 x_1 + 10^1 \cdot 0.1000 x_2 &= 10^1 \cdot 0.1000, \\ 10^1 \cdot 0.2000 x_1 + 10^1 \cdot 0.1000 x_2 &= 10^1 \cdot 0.2000. \end{aligned}$$

若用 $\frac{1}{2}(10^{-4} \cdot 0.1000)$ 除第一方程减第二方程, 则出现用小的数除大的数, 得到

$$\begin{aligned} 10^{-4} \cdot 0.1000 x_1 + 10^1 \cdot 0.1000 x_2 &= 10^1 \cdot 0.1000, \\ 10^6 \cdot 0.2000 x_2 &= 10^6 \cdot 0.2000. \end{aligned}$$

由此解出

$$x_1 = 0, \quad x_2 = 10^1 \cdot 0.1000 = 1,$$

显然严重失真.

若反过来用第二个方程消去第一个方程中含 x_1 的项, 则避免了大数被小数除, 得到

$$\begin{aligned} 10^6 \cdot 0.1000 x_2 &= 10^6 \cdot 0.1000, \\ 10^1 \cdot 0.2000 x_1 + 10^1 \cdot 0.1000 x_2 &= 10^1 \cdot 0.2000. \end{aligned}$$

由此求得相当好的近似解

$$x_1 = 0.5000, \quad x_2 = 10^1 \cdot 0.1000.$$

2. 要避免两相近数相减

在数值计算中两相近数相减有效数字会严重损失. 例如, $x = 532.65$, $y = 532.52$ 都具有五位有效数字, 但 $x - y = 0.13$ 只有两位有效数字. 这说明必须尽量避免出现这类运算. 最好是改变计算方法, 防止这种现象产生. 现举例说明.

例 7 求 $x^2 - 16x + 1 = 0$ 的小正根.

$$\text{解 } x_1 = 8 + \sqrt{63}, \quad x_2 = 8 - \sqrt{63} \quad 8 - \sqrt{7.94} = 0.06 = x_2^*,$$

x_2^* 只有一位有效数字 . 若改用

$$x_2 = 8 - 63 = \frac{1}{8 + 63} = \frac{1}{15.94} = 0.0627$$

具有 3 位有效数字 .

例 8 计算 $A = 10^7 (1 - \cos 2^\circ)$ (用四位数学用表) .

由于 $\cos 2^\circ = 0.9994$, 直接计算

$$A = 10^7 (1 - \cos 2^\circ) = 10^7 (1 - 0.9994) = 6 \times 10^3 .$$

只有一位有效数字 . 若利用 $1 - \cos x = 2 \sin^2 \frac{x}{2}$, 则

$$A = 10^7 (1 - \cos 2^\circ) = 2 \times (\sin 1^\circ)^2 \times 10^7 = 6.13 \times 10^3$$

具有三位有效数字(这里 $\sin 1^\circ = 0.0175$) .

此例说明, 可通过改变计算公式避免或减少有效数字的损失 .

类似地, 如果 x_1 和 x_2 很接近时, 则

$$\lg x_1 - \lg x_2 = \lg \frac{x_1}{x_2} .$$

用右边算式有效数字就不损失 . 当 x 很大时,

$$x + 1 - x = \frac{1}{x + 1 + x},$$

都用右端算式代替左端 . 一般情况, 当 $f(x) - f(x^*)$ 时, 可用泰勒展开

$$f(x) - f(x^*) = f(x^*)(x - x^*) + \frac{f'(x^*)}{2}(x - x^*)^2 + \dots$$

取右端的有限项近似左端 . 如果无法改变算式, 则采用增加有效位数进行运算; 在计算机上则采用双倍字长运算, 但这要增加机器计算时间和多占内存单元 .

3. 要防止大数“吃掉”小数

在数值运算中参加运算的数有时数量级相差很大, 而计算机位数有限, 如不注意运算次序就可能出现大数“吃掉”小数的现象, 影响计算结果的可靠性 .

例 9 在五位十进制计算机上, 计算

$$A = 52492 + \sum_{i=1}^{1000}$$

其中 $0.1 \dots_i 0.9$.

把运算的数写成规格化形式

$$A = 0.52492 \times 10^5 + \sum_{i=1}^{1000}$$

由于在计算机内计算时要对阶, 若取 $i = 0.9$, 对阶时 $i = 0.000009 \times 10^5$, 在五位的计算机中表示为机器 0, 因此

$$A = 0.52492 \times 10^5 + 0.000009 \times 10^5 + \dots + 0.000009 \times 10^5$$

$C 0.52492 \times 10^5$ (符号 C 表示机器中相等),

结果显然不可靠, 这是由于运算中出现了大数 52492“吃掉”小数造成的. 如果计算时先把数量级相同的一千个 i 相加, 最后再加 52492, 就不会出现大数“吃”小数现象, 这时

$$0.1 \times 10^3 + \sum_{i=1}^{1000} 0.9 \times 10^3,$$

于是

$$\begin{aligned} 0.001 \times 10^5 + 0.52492 \times 10^5 & A & 0.009 \times 10^5 + 0.52492 \times 10^5, \\ 52592 & A & 53392. \end{aligned}$$

4. 注意简化计算步骤, 减少运算次数

同样一个计算问题, 如果能减少运算次数, 不但可节省计算机的计算时间, 还能减少舍入误差. 这是数值计算必须遵从的原则, 也是“数值分析”要研究的重要内容.

例 10 计算多项式

$$P_n(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0$$

的值, 若直接计算 $a_k x^k$ 再逐项相加, 一共需做

$$n + (n - 1) + \dots + 2 + 1 = \frac{n(n+1)}{2}$$

次乘法和 n 次加法 . 若采用秦九韶算法

$$\begin{aligned} S_n &= a_n, \\ S_k &= xS_{k+1} + a_k \quad (k = n - 1, \dots, 2, 1, 0), \\ P_n(x) &= S_0. \end{aligned} \quad (3.4)$$

只要 n 次乘法和 n 次加法就可算出 $P_n(x)$ 的值 .

在“数值分析”中, 这种节省计算次数的算法还有不少 . 本书第 3 章介绍的 FFT 算法, 就是一个最成功的范例 .

评 注

本章第 1 节简要地介绍了数值分析的研究对象, 它是计算数学的主要部分, 关于计算数学介绍可参考《中国大百科全书·数学》中的有关条目 . 第 2 和第 3 节中介绍了误差的基本概念与误差分析的若干原则, 数值计算中的舍入误差是一个困难而复杂的问题, 目前尚无真正有效的定量估计方法, 第 3 节中提到的向后误差分析法(见文献[8])及区间分析法^[9], 实际计算时仍有不少困难, 对大型科学计算的误差估计仍难于使用 . 因此在本书中更着重对误差的定性分析, 即对每个具体算法只要是数值稳定的, 就不必再做舍入误差估计, 至于方法的截断误差将结合不同问题的具体算法进行讨论 .

习 题

- 1 . 设 $x > 0$, x 的相对误差为 , 求 $\ln x$ 的误差 .
- 2 . 设 x 的相对误差为 2 %, 求 x^n 的相对误差 .
- 3 . 下列各数都是经过四舍五入得到的近似数, 即误差限不超过最后一位的半个单位, 试指出它们是几位有效数字:

$$\begin{aligned}x_1^* &= 1.1021, & x_2^* &= 0.031, & x_3^* &= 385.6, \\x_4^* &= 56.430, & x_5^* &= 7 \times 1.0.\end{aligned}$$

4. 利用公式(2.3)求下列各近似值的误差限:

$$(i) x_1^* + x_2^* + x_4^*, \quad (ii) x_1^* x_2^* x_3^*, \quad (iii) x_2^* / x_4^*.$$

其中 x_1^* , x_2^* , x_3^* , x_4^* 均为第 3 题所给的数.

5. 计算球体积要使相对误差限为 1%, 问度量半径 R 时允许的相对误差限是多少?

6. 设 $Y_0 = 28$, 按递推公式

$$Y_n = Y_{n-1} - \frac{1}{100} \cdot 783 \quad (n = 1, 2, \dots)$$

计算到 Y_{100} . 若取 783 27.982(5 位有效数字), 试问计算 Y_{100} 将有多大误差?

7. 求方程 $x^2 - 56x + 1 = 0$ 的两个根, 使它至少具有 4 位有效数字
(783 27.982).

8. 当 N 充分大时, 怎样求 $\int_N^{N+1} \frac{1}{1+x^2} dx$?

9. 正方形的边长大约为 100cm, 应怎样测量才能使其面积误差不超过 1cm^2 ?

10. 设 $S = \frac{1}{2}gt^2$, 假定 g 是准确的, 而对 t 的测量有 ± 0.1 秒的误差, 证明当 t 增加时 S 的绝对误差增加, 而相对误差却减少.

11. 序列 $\{y_n\}$ 满足递推关系

$$y_n = 10y_{n-1} - 1 \quad (n = 1, 2, \dots),$$

若 $y_0 = 2.1.41$ (三位有效数字), 计算到 y_{10} 时误差有多大? 这个计算过程稳定吗?

12. 计算 $f = (2 - 1)^6$, 取 2 1.4, 利用下列等式计算, 哪一个得到的结果最好?

$$\frac{1}{(2+1)^6}, \quad (3-2-2)^3,$$

$$\frac{1}{(3+2-2)^3}, \quad 99 - 70 - 2.$$

13. $f(x) = \ln(x - \sqrt{x^2 - 1})$, 求 $f(30)$ 的值. 若开平方用 6 位函数表, 问求对数时误差有多大? 若改用另一等价公式

$$\ln(x - \sqrt{x^2 - 1}) = -\ln(x + \sqrt{x^2 - 1})$$

计算, 求对数时误差有多大?

第2章 插 值 法

2.1 引 言

许多实际问题都用函数 $y = f(x)$ 来表示某种内在规律的数量关系, 其中相当一部分函数是通过实验或观测得到的. 虽然 $f(x)$ 在某个区间 $[a, b]$ 上是存在的, 有的还是连续的, 但却只能给出 $[a, b]$ 上一系列点 x_i 的函数值 $y_i = f(x_i)$ ($i = 0, 1, \dots, n$), 这只是一张函数表. 有的函数虽有解析表达式, 但由于计算复杂, 使用不方便, 通常也造一个函数表, 如大家熟悉的三角函数表、对数表、平方根和立方根表等等. 为了研究函数的变化规律, 往往需要求出不在表上的函数值. 因此, 我们希望根据给定的函数表做一个既能反映函数 $f(x)$ 的特性, 又便于计算的简单函数 $P(x)$, 用 $P(x)$ 近似 $f(x)$. 通常选一类较简单的函数(如代数多项式或分段代数多项式)作为 $P(x)$, 并使 $P(x_i) = f(x_i)$ 对 $i = 0, 1, \dots, n$ 成立. 这样确定的 $P(x)$ 就是我们希望得到的插值函数. 例如, 在现代机械工业中用计算机程序控制加工机械零件, 根据设计可给出零件外形曲线的某些型值点 (x_i, y_i) ($i = 0, 1, \dots, n$), 加工时为控制每步走刀方向及步数, 就要算出零件外形曲线其他点的函数值, 才能加工出外表光滑的零件, 这就是求插值函数的问题. 下面我们给出有关插值法的定义.

设函数 $y = f(x)$ 在区间 $[a, b]$ 上有定义, 且已知在点 $a = x_0 < x_1 < \dots < x_n = b$ 上的值 y_0, y_1, \dots, y_n , 若存在一简单函数 $P(x)$, 使

$$P(x_i) = y_i \quad (i = 0, 1, \dots, n) \quad (1.1)$$

成立, 就称 $P(x)$ 为 $f(x)$ 的插值函数, 点 x_0, x_1, \dots, x_n 称为插值节

点,包含插值节点的区间 $[a, b]$ 称为插值区间,求插值函数 $P(x)$ 的方法称为插值法.若 $P(x)$ 是次数不超过 n 的代数多项式,即

$$P(x) = a_0 + a_1 x + \dots + a_n x^n, \quad (1.2)$$

其中 a_i 为实数,就称 $P(x)$ 为插值多项式,相应的插值法称为多项式插值.若 $P(x)$ 为分段的多项式,就称为分段插值.若 $P(x)$ 为三角多项式,就称为三角插值.本章只讨论多项式插值与分段插值.

从几何上看,插值法就是求曲线 $y = P(x)$,使其通过给定的 $n+1$ 个点 (x_i, y_i) , $i=0, 1, \dots, n$,并用它近似已知曲线 $y = f(x)$,见图2-1.

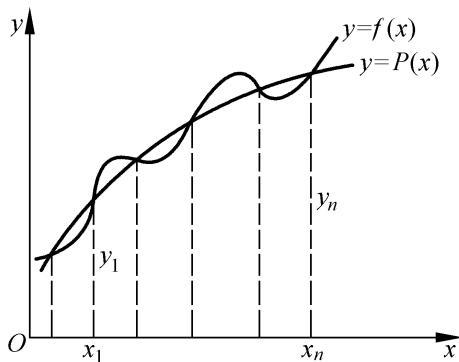


图 2-1

插值法是一种古老的数学方法,它来自生产实践.早在一千多年前,我国科学家在研究历法上就应用了线性插值与二次插值,但它的基本理论和结果却是在微积分产生以后才逐步完善的,其应用也日益增多,特别是在电子计算机广泛使用以后,由于航空、造船、精密机械加工等实际问题的需要,使插值法在实践上或理论上显得更为重要,并得到进一步发展,尤其是近几十年发展起来的样条(spline)插值,更获得了广泛的应用.

本章主要研究如何求出插值多项式,分段插值函数,样条插值函数;讨论插值多项式 $P(x)$ 的存在唯一性、收敛性及误差估计等.

2.2 拉格朗日插值

2.2.1 线性插值与抛物插值

对给定的插值点为求得形如(1.2)的插值多项式可以有各种不同方法,下面先讨论 $n=1$ 的简单情形,假定给定区间 $[x_k, x_{k+1}]$ 及端点函数值 $y_k = f(x_k)$, $y_{k+1} = f(x_{k+1})$, 要求线性插值多项式 $L(x)$, 使它满足

$$L_1(x_k) = y_k, \quad L_1(x_{k+1}) = y_{k+1}.$$

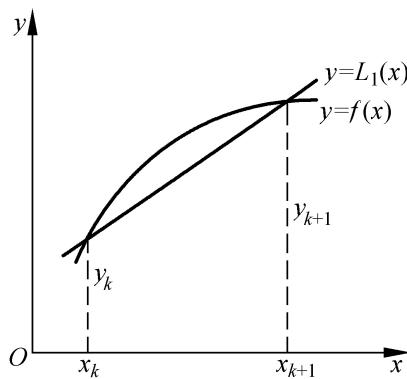


图 2-2

$y = L_1(x)$ 的几何意义就是通过两点 (x_k, y_k) 与 (x_{k+1}, y_{k+1}) 的直线, 如图 2-2 所示, $L_1(x)$ 的表达式可由几何意义直接给出

$$\begin{aligned} L_1(x) &= y_k + \frac{y_{k+1} - y_k}{x_{k+1} - x_k} (x - x_k) \quad (\text{点斜式}), \\ L_1(x) &= \frac{x_{k+1} - x}{x_{k+1} - x_k} y_k + \frac{x - x_k}{x_{k+1} - x_k} y_{k+1} \quad (\text{两点式}). \end{aligned} \quad (2.1)$$

由两点式看出, $L_1(x)$ 是由两个线性函数

$$l_k(x) = \frac{x - x_{k+1}}{x_k - x_{k+1}}, \quad l_{k+1}(x) = \frac{x - x_k}{x_{k+1} - x_k} \quad (2.2)$$

的线性组合得到, 其系数分别为 y_k 及 y_{k+1} , 即

$$L(x) = y_k l_k(x) + y_{k+1} l_{k+1}(x). \quad (2.3)$$

显然, $l_k(x)$ 及 $l_{k+1}(x)$ 也是线性插值多项式, 在节点 x_k 及 x_{k+1} 上满足条件

$$\begin{aligned} l_k(x_k) &= 1, & l_k(x_{k+1}) &= 0; \\ l_{k+1}(x_k) &= 0, & l_{k+1}(x_{k+1}) &= 1. \end{aligned}$$

我们称函数 $l_k(x)$ 及 $l_{k+1}(x)$ 为线性插值基函数, 它们的图形见图 2-3.

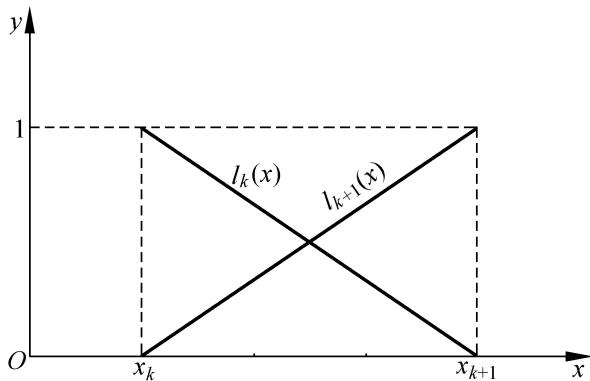


图 2-3

下面讨论 $n=2$ 的情况. 假定插值节点为 x_{k-1} , x_k , x_{k+1} , 要求二次插值多项式 $L_2(x)$, 使它满足

$$L_2(x_j) = y_j \quad (j = k-1, k, k+1).$$

我们知道 $y = L_2(x)$ 在几何上就是通过三点 (x_{k-1}, y_{k-1}) , (x_k, y_k) , (x_{k+1}, y_{k+1}) 的抛物线. 为了求出 $L_2(x)$ 的表达式, 可采用基函数方法, 此时基函数 $l_{k-1}(x)$, $l_k(x)$ 及 $l_{k+1}(x)$ 是二次函数, 且在节点上满足条件

$$\begin{aligned} l_{k-1}(x_{k-1}) &= 1, & l_{k-1}(x_j) &= 0 \quad (j = k, k+1); \\ l_k(x_k) &= 1, & l_k(x_j) &= 0 \quad (j = k-1, k+1); \\ l_{k+1}(x_{k+1}) &= 1, & l_{k+1}(x_j) &= 0 \quad (j = k-1, k). \end{aligned} \quad (2.4)$$

满足条件(2.4)的插值基函数是很容易求出的, 例如求 $l_{k-1}(x)$, 因它有两个零点 x_k 及 x_{k+1} , 故可表示为

$$l_{k-1}(x) = A(x - x_k)(x - x_{k+1}),$$

其中 A 为待定系数, 可由条件 $l_{k-1}(x_{k-1}) = 1$ 定出

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

于是 $l_{k-1}(x) = \frac{(x - x_k)(x - x_{k+1})}{(x_{k-1} - x_k)(x_{k-1} - x_{k+1})}.$

同理可得 $l_k(x) = \frac{(x - x_{k-1})(x - x_{k+1})}{(x_k - x_{k-1})(x_k - x_{k+1})},$

$$l_{k+1}(x) = \frac{(x - x_{k-1})(x - x_k)}{(x_{k+1} - x_{k-1})(x_{k+1} - x_k)}.$$

二次插值基函数 $l_{k-1}(x)$, $l_k(x)$, $l_{k+1}(x)$ 在区间 $[x_{k-1}, x_{k+1}]$ 上的图形见图 2-4.

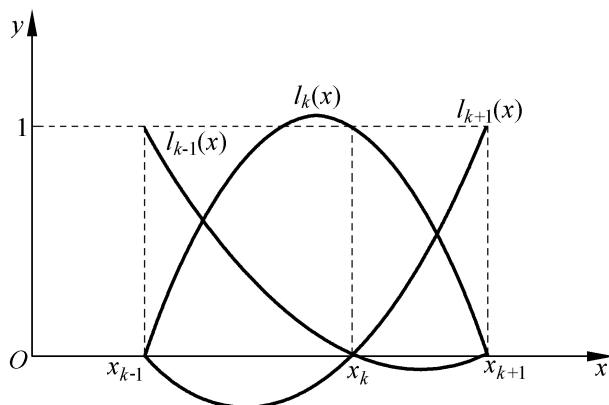


图 2-4

利用二次插值基函数 $l_{k-1}(x)$, $l_k(x)$, $l_{k+1}(x)$, 立即得到二次插值多项式

$$L_2(x) = y_{k-1} l_{k-1}(x) + y_k l_k(x) + y_{k+1} l_{k+1}(x), \quad (2.5)$$

显然, 它满足条件 $L_2(x_j) = y_j$ ($j = k-1, k, k+1$). 将上面求得的 $l_{k-1}(x)$, $l_k(x)$, $l_{k+1}(x)$ 代入 (2.5), 得

$$L_2(x) = y_{k-1} \frac{(x - x_k)(x - x_{k+1})}{(x_{k-1} - x_k)(x_{k-1} - x_{k+1})}$$

$$\begin{aligned}
 & + y_k \frac{(x - x_{k-1})(x - x_{k+1})}{(x_k - x_{k-1})(x_k - x_{k+1})} \\
 & + y_{k+1} \frac{(x - x_{k-1})(x - x_k)}{(x_{k+1} - x_{k-1})(x_{k+1} - x_k)} .
 \end{aligned}$$

2.2.2 拉格朗日插值多项式

上面我们对 $n=1$ 及 $n=2$ 的情况, 得到了一次与二次插值多项式 $L_1(x)$ 及 $L_2(x)$, 它们分别由(2.3)与(2.5)表示. 这种用插值基函数表示的方法容易推广到一般情形. 下面讨论如何构造通过 $n+1$ 个节点 $x_0 < x_1 < \dots < x_n$ 的 n 次插值多项式 $L_n(x)$, 假定它满足条件

$$L_n(x_j) = y_j \quad (j = 0, 1, \dots, n). \quad (2.6)$$

为了构造 $L_n(x)$, 我们先定义 n 次插值基函数.

定义 1 若 n 次多项式 $l_j(x) (j = 0, 1, \dots, n)$ 在 $n+1$ 个节点 $x_0 < x_1 < \dots < x_n$ 上满足条件

$$l_j(x_k) = \begin{cases} 1, & k = j; \\ 0, & k \neq j. \end{cases} \quad (j, k = 0, 1, \dots, n) \quad (2.7)$$

就称这 $n+1$ 个 n 次多项式 $l_0(x), l_1(x), \dots, l_n(x)$ 为节点 x_0, x_1, \dots, x_n 上的 n 次插值基函数.

当 $n=1$ 及 $n=2$ 时的情况前面已经讨论. 用类似的推导方法, 可得到 n 次插值基函数为

$$\begin{aligned}
 l_k(x) = & \frac{(x - x_0) \dots (x - x_{k-1})(x - x_{k+1}) \dots (x - x_n)}{(x_k - x_0) \dots (x_k - x_{k-1})(x_k - x_{k+1}) \dots (x_k - x_n)} \\
 & (k = 0, 1, \dots, n) . \quad (2.8)
 \end{aligned}$$

显然它满足条件(2.7). 于是, 满足条件(2.6)的插值多项式 $L_n(x)$ 可表示为

$$L_n(x) = \sum_{k=0}^n y_k l_k(x) . \quad (2.9)$$

由 $l_k(x)$ 的定义, 知

$$L_n(x_j) = \sum_{k=0}^n y_k l_k(x_j) = y_j \quad (j = 0, 1, \dots, n).$$

形如(2.9)的插值多项式 $L_n(x)$ 称为拉格朗日(Lagrange)插值多项式, 而(2.3)与(2.5)是 $n=1$ 和 $n=2$ 的特殊情形.

若引入记号

$$n+1(x) = (x - x_0)(x - x_1) \dots (x - x_n), \quad (2.10)$$

容易求得

$$n+1(x_k) = (x_k - x_0) \dots (x_k - x_{k-1})(x_k - x_{k+1}) \dots (x_k - x_n).$$

于是公式(2.9)可改写成

$$L_n(x) = \sum_{k=0}^n y_k \frac{n+1(x)}{(x - x_k)_{n+1}(x_k)}. \quad (2.11)$$

注意, n 次插值多项式 $L_n(x)$ 通常是次数为 n 的多项式, 特殊情况下次数可能小于 n . 例如, 通过三点 $(x_0, y_0), (x_1, y_1), (x_2, y_2)$ 的二次插值多项式 $L_2(x)$, 如果三点共线, 则 $y = L_2(x)$ 就是一直线, 而不是抛物线, 这时 $L_2(x)$ 是一次多项式.

关于插值多项式存在唯一性有以下定理.

定理 1 在次数不超过 n 的多项式集合 H_n 中, 满足条件(2.6)的插值多项式 $L_n(x) \in H_n$ 是存在唯一的.

证明 公式(2.11)所表示的 $L_n(x)$ 已证明了插值多项式的存在性, 下面用反证法证明唯一性. 假定还有 $P(x) \in H_n$ 使 $P(x_i) = f(x_i), i = 0, 1, \dots, n$ 成立. 于是有 $L_n(x_i) - P(x_i) = 0$ 对 $i = 0, 1, \dots, n$ 成立, 它表明多项式 $L_n(x) - P(x) \in H_n$ 有 $n+1$ 个零点 x_0, x_1, \dots, x_n 这与 n 次多项式只有 n 个零点的代数基本定理矛盾, 故只能 $P(x) = L_n(x)$. 证毕.

根据存在唯一性定理, 若令 $f(x) = x^m, m = 0, 1, \dots, n$ 可得

$$\sum_{k=0}^n x_k^m l_k(x) = x^m, \quad m = 0, 1, \dots, n. \quad (2.12)$$

若取 $m=0$, 则

$$\sum_{k=0}^n l_k(x) = 1. \quad (2.13)$$

它可用来检验函数组 $\{l_k(x), k=0, 1, \dots, n\}$ 的正确性.

2.2.3 插值余项与误差估计

若在 $[a, b]$ 上用 $L_n(x)$ 近似 $f(x)$, 则其截断误差为 $R_n(x) = f(x) - L_n(x)$, 也称为插值多项式的余项. 关于插值余项估计有以下定理.

定理 2 设 $f^{(n)}(x)$ 在 $[a, b]$ 上连续, $f^{(n+1)}(x)$ 在 (a, b) 内存在, 节点 $a = x_0 < x_1 < \dots < x_n = b$, $L_n(x)$ 是满足条件 (2.6) 的插值多项式, 则对任何 $x \in [a, b]$, 插值余项

$$R_n(x) = f(x) - L_n(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} (x - x_0)(x - x_1)\dots(x - x_n), \quad (2.14)$$

这里 $\xi \in (a, b)$ 且依赖于 x , $\xi^{(n+1)}(x)$ 是 (2.10) 所定义的.

证明 由给定条件知 $R_n(x)$ 在节点 x_k ($k=0, 1, \dots, n$) 上为零, 即 $R_n(x_k) = 0$ ($k=0, 1, \dots, n$), 于是

$$R_n(x) = K(x)(x - x_0)(x - x_1)\dots(x - x_n) = K(x)\xi^{(n+1)}(x), \quad (2.15)$$

其中 $K(x)$ 是与 x 有关的待定函数.

现把 x 看成 $[a, b]$ 上的一个固定点, 作函数

$$(t) = f(t) - L_n(t) - K(x)(t - x_0)(t - x_1)\dots(t - x_n),$$

根据插值条件及余项定义, 可知 (t) 在点 x_0, x_1, \dots, x_n 及 x 处均为零, 故 (t) 在 $[a, b]$ 上有 $n+2$ 个零点, 根据罗尔 (Rolle) 定理,

(t) 在 (t) 的两个零点间至少有一个零点, 故 (t) 在 $[a, b]$ 内至少有 $n+1$ 个零点. 对 (t) 再应用罗尔定理, 可知 (t) 在 $[a, b]$ 内至少有 n 个零点. 依此类推, $\xi^{(n+1)}(t)$ 在 (a, b) 内至少有一个零点, 记为 (a, b) , 使

$$\xi^{(n+1)}(\xi) = f^{(n+1)}(\xi) - (n+1)!K(x) = 0,$$

于是

$$K(x) = \frac{f^{(n+1)}(\cdot)}{(n+1)!}, \quad (a, b), \text{且依赖于 } x.$$

将它代入(2.15), 就得到余项表达式(2.14). 证毕.

应当指出, 余项表达式只有在 $f(x)$ 的高阶导数存在时才能应用. 在 (a, b) 内的具体位置通常不可能给出, 如果我们可以求出 $\max_{a < x < b} |f^{(n+1)}(x)| = M_{n+1}$, 那么插值多项式 $L_n(x)$ 逼近 $f(x)$ 的截断误差限是

$$|R_n(x)| \leq \frac{M_{n+1}}{(n+1)!} |L_{n+1}(x)|. \quad (2.16)$$

当 $n=1$ 时, 线性插值余项为

$$R_1(x) = \frac{1}{2} f''(\cdot)(x - x_0)(x - x_1),$$

$$[x_0, x_1]; \quad (2.17)$$

当 $n=2$ 时, 抛物插值的余项为

$$R_2(x) = \frac{1}{6} f'''(\cdot)(x - x_0)(x - x_1)(x - x_2),$$

$$[x_0, x_2]. \quad (2.18)$$

例 1 已给 $\sin 0.32 = 0.314567$, $\sin 0.34 = 0.333487$, $\sin 0.36 = 0.352274$, 用线性插值及抛物插值计算 $\sin 0.3367$ 的值并估计截断误差.

解 由题意取 $x_0 = 0.32$, $y_0 = 0.314567$, $x_1 = 0.34$, $y_1 = 0.333487$, $x_2 = 0.36$, $y_2 = 0.352274$.

用线性插值计算, 取 $x_0 = 0.32$ 及 $x_1 = 0.34$, 由公式(2.1)得

$$\begin{aligned} \sin 0.3367 - L(0.3367) &= y_0 + \frac{y_1 - y_0}{x_1 - x_0}(0.3367 - x_0) \\ &= 0.314567 + \frac{0.01892}{0.02} \times 0.0167 = 0.330365. \end{aligned}$$

其截断误差由(2.17)得

$$| R_1(x) | \leq \frac{M_2}{2} / (x - x_0)(x - x_1),$$

其中 $M_2 = \max_{x_0 \leq x \leq x_1} |f(x)|$, 因 $f(x) = \sin x$, $f'(x) = -\cos x$. 可取 $M_2 = \max_{x_0 \leq x \leq x_1} |\sin x| = \sin x_1 = 0.3335$, 于是

$$\begin{aligned} |R_1(0.3367)| &= |\sin 0.3367 - L_1(0.3367)| \\ &\leq \frac{1}{2} \times 0.3335 \times 0.0167 \times 0.0033 = 0.92 \times 10^{-5}. \end{aligned}$$

用抛物插值计算 $\sin 0.3367$ 时, 由公式(2.5)得

$$\begin{aligned} \sin 0.3367 &= y_0 \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} + y_1 \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)} \\ &\quad + y_2 \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)} \\ &= L_2(0.3367) \\ &= 0.314567 \times \frac{0.7689 \times 10^{-4}}{0.0008} + 0.333487 \\ &\quad \times \frac{3.89 \times 10^{-4}}{0.0004} + 0.352274 \times \frac{-0.5511 \times 10^{-4}}{0.0008} \\ &= 0.330374. \end{aligned}$$

这个结果与 6 位有效数字的正弦函数表完全一样, 这说明查表时用二次插值精度已相当高了. 其截断误差限由(2.18)得

$$|R_2(x)| \leq \frac{M_3}{6} / (x - x_0)(x - x_1)(x - x_2),$$

其中

$$M_3 = \max_{x_0 \leq x \leq x_2} |f'(x)| = \cos x_0 < 0.828,$$

于是

$$\begin{aligned} |R_2(0.3367)| &= |\sin 0.3367 - L_2(0.3367)| \\ &\leq \frac{1}{6} \times 0.828 \times 0.0167 \times 0.033 \times 0.0233 \\ &< 0.178 \times 10^{-6}. \end{aligned}$$

2.3 均差与牛顿插值公式

2.3.1 均差及其性质

利用插值基函数很容易得到拉格朗日插值多项式, 公式结构紧凑, 在理论分析中甚为方便, 但当插值节点增减时全部插值基函数 $l_k(x) (k=0, 1, \dots, n)$ 均要随之变化, 整个公式也将发生变化, 这在实际计算中是很不方便的, 为了克服这一缺点, 可把插值多项式表示为如下便于计算的形式

$$\begin{aligned} P_n(x) = & a_0 + a_1(x - x_0) + a_2(x - x_0)(x - x_1) + \dots \\ & + a_n(x - x_0) \dots (x - x_{n-1}), \end{aligned} \quad (3.1)$$

其中 a_0, a_1, \dots, a_n 为待定系数, 可由插值条件

$$P_n(x_j) = f_j \quad (j = 0, 1, \dots, n)$$

确定.

当 $x = x_0$ 时, $P_n(x_0) = a_0 = f_0$.

当 $x = x_1$ 时, $P_n(x_1) = a_0 + a_1(x_1 - x_0) = f_1$, 推得

$$a_1 = \frac{f_1 - f_0}{x_1 - x_0}.$$

当 $x = x_2$ 时, $P_n(x_2) = a_0 + a_1(x_2 - x_0) + a_2(x_2 - x_0)(x_2 - x_1) = f_2$, 推得

$$a_2 = \frac{\frac{f_2 - f_0}{x_2 - x_0} - \frac{f_1 - f_0}{x_1 - x_0}}{x_2 - x_1}.$$

依此递推可得到 a_0, \dots, a_n . 为写出系数 a_k 的一般表达式, 先引进如下均差定义.

定义 2 称 $f[x_0, x_k] = \frac{f(x_k) - f(x_0)}{x_k - x_0}$ 为函数 $f(x)$ 关于点

x_0, x_k 的一阶均差. $f[x_0, x_1, x_k] = \frac{f[x_0, x_k] - f[x_0, x_1]}{x_k - x_1}$ 称为 $f(x)$ 的二阶均差. 一般地, 称

$$f[x_0, x_1, \dots, x_k] = \frac{f[x_0, \dots, x_{k-2}, x_k] - f[x_0, \dots, x_{k-1}]}{x_k - x_{k-1}} \quad (3.2)$$

为 $f(x)$ 的 k 阶均差(均差也称为差商).

均差有如下的基本性质:

1° k 阶均差可表为函数值 $f(x_0), \dots, f(x_k)$ 的线性组合, 即

$$f[x_0, \dots, x_k] = \frac{f(x_j)}{(x_j - x_0) \dots (x_j - x_{j-1})(x_j - x_{j+1}) \dots (x_j - x_k)}. \quad (3.3)$$

这个性质可用归纳法证明. 这性质也表明均差与节点的排列次序无关, 称为均差的对称性. 即

$$f[x_0, \dots, x_k] = f[x_1, x_0, x_2, \dots, x_k] = \dots = f[x_1, \dots, x_k, x_0].$$

2° 由性质 1° 及 (3.2) 可得

$$f[x_0, \dots, x_k] = \frac{f[x_1, \dots, x_k] - f[x_0, \dots, x_{k-1}]}{x_k - x_0}. \quad (3.4)$$

3° 若 $f(x)$ 在 $[a, b]$ 上存在 n 阶导数, 且节点 x_0, \dots, x_n $[a, b]$, 则 n 阶均差与导数关系如下:

$$f[x_0, \dots, x_n] = \frac{f^{(n)}(\cdot)}{n!}, \quad [a, b]. \quad (3.5)$$

这公式可直接用罗尔定理证明.

均差的其他性质还可见习题. 均差计算可列均差表如下
(表 2-1).

表 2.1

x_k	$f(x_k)$	一阶均差	二阶均差	三阶均差	四阶均差
x_0	$f(x_0)$				
x_1	$f(x_1)$	$f[x_0, x_1]$			
x_2	$f(x_2)$	$f[x_1, x_2]$	$f[x_0, x_1, x_2]$		
x_3	$f(x_3)$	$f[x_2, x_3]$	$f[x_0, x_1, x_2, x_3]$	$f[x_0, x_1, x_2, x_3]$	
x_4	$f(x_4)$	$f[x_3, x_4]$	$f[x_2, x_3, x_4]$	$f[x_1, x_2, x_3, x_4]$	$f[x_0, x_1, x_2, x_3, x_4]$
...

2.3.2 牛顿插值公式

根据均差定义, 把 x 看成 $[a, b]$ 上一点, 可得

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

...

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

只要把后一式代入前一式, 就得到

$$\begin{aligned} f(x) &= f(x_0) + f[x_0, x_1](x - x_0) \\ &\quad + f[x_0, x_1, x_2](x - x_0)(x - x_1) + \dots \\ &\quad + f[x_0, x_1, \dots, x_n](x - x_0) \dots (x - x_{n-1}) \\ &\quad + f[x, x_0, \dots, x_n]_{n+1}(x) = N_n(x) + R_n(x), \end{aligned}$$

其中

$$\begin{aligned} N_n(x) &= f(x_0) + f[x_0, x_1](x - x_0) \\ &\quad + f[x_0, x_1, x_2](x - x_0)(x - x_1) + \dots \\ &\quad + f[x_0, \dots, x_n](x - x_0) \dots (x - x_{n-1}), \quad (3.6) \end{aligned}$$

$$R_n(x) = f(x) - N_n(x) = f[x, x_0, \dots, x_n]_{n+1}(x), \quad (3.7)$$

$n+1$ 次多项式是由(2.10)定义的.

由(3.6)确定的多项式 $N_n(x)$ 显然满足插值条件, 且次数不

超过 n , 它就是形如(3.1)的多项式, 其系数为

$$a_k = f[x_0, \dots, x_k] \quad (k = 0, 1, \dots, n).$$

我们称 $N_n(x)$ 为牛顿(Newton)均差插值多项式. 系数 a_k 就是均差表 2-1 中加横线的各阶均差, 它比拉格朗日插值计算量省, 且便于程序设计.

(3.7) 为插值余项, 由插值多项式唯一性知, 它与(2.14)是等价的, 事实上, 利用均差与导数关系式(3.5)可由(3.7)推出(2.14). 但(3.7)更有一般性, 它对 f 是由离散点给出的情形或 f 导数不存在时均适用.

例 2 给出 $f(x)$ 的函数表(见表 2-2), 求 4 次牛顿插值多项式, 并由此计算 $f(0.596)$ 的近似值.

首先根据给定函数表造出均差表.

表 2-2

0.40	0.41075					
0.55	0.57815	1.11600				
0.65	0.69675	1.18600	0.28000			
0.80	0.88811	1.27573	0.35893	0.19733		
0.90	1.02652	1.38410	0.43348	0.21300	0.03134	
1.05	1.25382	1.51533	0.52493	0.22863	0.03126	-0.00012

从均差表看到 4 阶均差近似常数. 故取 4 次插值多项式 $N_4(x)$ 做近似即可.

$$\begin{aligned} N_4(x) &= 0.41075 + 1.116(x - 0.4) + 0.28(x - 0.4)(x - 0.55) \\ &\quad + 0.19733(x - 0.4)(x - 0.55)(x - 0.65) \\ &\quad + 0.03134(x - 0.4)(x - 0.55)(x - 0.65)(x - 0.8), \end{aligned}$$

于是

$$f(0.596) \quad N_4(0.596) = 0.63192,$$

截断误差

$$|R_4(x)| / |f[x_0, \dots, x_5](0.596)| = 3.63 \times 10^{-9}.$$

这说明截断误差很小, 可忽略不计.

此例的截断误差估计中, 5 阶均差 $f[x, x_0, \dots, x_4]$ 用 $f[x_0, x_1, \dots, x_5] = -0.00012$ 近似. 另一种方法是取 $x = 0.596$, 由 $f(0.596) = 0.63192$, 可求得 $f[x, x_0, \dots, x_4]$ 的近似值, 从而可求得 $|R_4(x)|$ 的近似.

2.4 差分与等距节点插值

上面讨论了节点任意分布的插值公式, 但实际应用时经常遇到等距节点的情形, 这时插值公式可以进一步简化, 计算也简单得多. 为了得到等距节点的插值公式, 我们先介绍差分的概念.

2.4.1 差分及其性质

设函数 $y = f(x)$ 在等距节点 $x_k = x_0 + kh (k=0, 1, \dots, n)$ 上的值 $f_k = f(x_k)$ 为已知, 这里 h 为常数, 称为步长.

定义 3 记号

$$f_k = f_{k+1} - f_k, \quad (4.1)$$

$$f_k = f_k - f_{k-1}, \quad (4.2)$$

$$f_k = f(x_k + h/2) - f(x_k - h/2) = f_{k+\frac{1}{2}} - f_{k-\frac{1}{2}} \quad (4.3)$$

分别称为 $f(x)$ 在 x_k 处以 h 为步长的向前差分, 向后差分及中心差分. 符号 Δ , ∇ , $\bar{\Delta}$ 分别称为向前差分算子, 向后差分算子及中心差分算子.

利用一阶差分可定义二阶差分为

$$\Delta^2 f_k = f_{k+1} - f_k = f_{k+2} - 2f_{k+1} + f_k.$$

一般地可定义 m 阶差分为

$$\Delta^m f_k = f_{k+1} - f_k;$$

$$\Delta^m f_k = f_k - f_{k-1}.$$

因中心差分 f_k 用到 $f_{k+\frac{1}{2}}$ 及 $f_{k-\frac{1}{2}}$ 这两个值, 实际上不是函数表上的值, 如果用函数表上的值, 一阶中心差分应写成

$$f_{k+\frac{1}{2}} = f_{k+1} - f_k, \quad f_{k-\frac{1}{2}} = f_k - f_{k-1},$$

二阶中心差分为

$$^2 f_k = f_{k+\frac{1}{2}} - f_{k-\frac{1}{2}},$$

等等.

除了已引入的差分算子外, 常用算子符号还有不变算子 I 及移位算子 E , 定义如下:

$$If_k = f_k, \quad Ef_k = f_{k+1},$$

于是, 由 $f_k = f_{k+1} - f_k = Ef_k - If_k = (E - I)f_k$, 可得

$$= E - I,$$

同理可得

$$= I - E^{-1}, \quad = E^{\frac{1}{2}} - E^{-\frac{1}{2}}.$$

由差分定义并应用算子符号运算可得下列基本性质.

性质 1 各阶差分均可用函数值表示. 例如

$$\begin{aligned} {}^n f_k &= (E - I)^n f_k = \sum_{j=0}^n (-1)^j \binom{n}{j} E^{n-j} f_k \\ &= \sum_{j=0}^n (-1)^j \binom{n}{j} f_{n+k-j}, \end{aligned} \quad (4.4)$$

$$\begin{aligned} {}^n f_k &= (I - E^{-1})^n f_k = \sum_{j=0}^n (-1)^{n-j} \binom{n}{j} E^{j-n} f_k \\ &= \sum_{j=0}^n (-1)^{n-j} \binom{n}{j} f_{k+j-n}, \end{aligned} \quad (4.5)$$

其中 $\binom{n}{j} = \frac{n(n-1)\dots(n-j+1)}{j!}$ 为二项式展开系数.

性质 2 可用各阶差分表示函数值. 例如, 可用向前差分表示 f_{n+k} , 因为

$$f_{n+k} = E^n f_k = (I +)^n f_k = \sum_{j=0}^n \binom{n}{j} f_k,$$

于是

$$f_{n+k} = \sum_{j=0}^n \binom{n}{j} f_k. \quad (4.6)$$

性质 3 均差与差分有密切关系, 例如, 对向前差分, 由定义

$$\begin{aligned}
 f[x_k, x_{k+1}] &= \frac{f_{k+1} - f_k}{x_{k+1} - x_k} = \frac{f_k}{h}, \\
 f[x_k, x_{k+1}, x_{k+2}] &= \frac{f[x_{k+1}, x_{k+2}] - f[x_k, x_{k+1}]}{x_{k+2} - x_k} \\
 &= \frac{1}{2h^2} f_k,
 \end{aligned}$$

一般地有

$$f[x_k, \dots, x_{k+m}] = \frac{1}{m!} \frac{1}{h^m} {}^m f_k \quad (m = 1, 2, \dots, n). \quad (4.7)$$

同理, 对向后差分有

$$f[x_k, x_{k-1}, \dots, x_{k-m}] = \frac{1}{m!} \frac{1}{h^m} {}^m f_k. \quad (4.8)$$

利用(4.7)及(3.5)又可得到

$${}^n f_k = h^n f^{(n)}(), \quad (4.9)$$

其中 (x_k, x_{k+n}) , 这就是差分与导数的关系. 差分的其他性质可参看本章习题.

计算差分可列差分表(见表 2-3), 表中 为向前差分, 为向后差分.

表 2-3

f_k	$()$	${}^2 ()$	${}^3 ()$	${}^4 ()$...
f_0	$f_0 (f_1)$				
f_1	$f_1 (f_2)$	${}^2 f_0 ({}^2 f_2)$	${}^3 f_0 ({}^3 f_3)$	${}^4 f_0 ({}^4 f_4)$	
f_2	$f_2 (f_3)$	${}^2 f_1 ({}^2 f_3)$	${}^3 f_1 ({}^3 f_4)$	${}^4 f_0 ({}^4 f_4)$...
f_3	$f_3 (f_4)$	${}^2 f_2 ({}^2 f_4)$	$... \dots$...	
f_4			
...					

2.4.2 等距节点插值公式

将牛顿均差插值多项式(3.6)中各阶均差用相应差分代替, 就可得到各种形式的等距节点插值公式. 这里只推导常用的前插与后插公式.

如果节点 $x_k = x_0 + kh$ ($k=0, 1, \dots, n$), 要计算 x_0 附近点 x 的函数 $f(x)$ 的值, 可令 $x = x_0 + th$, $0 \leq t \leq 1$, 于是

$$N_{k+1}(x) = \sum_{j=0}^k (x - x_j) = t(t-1)\dots(t-k) h^{k+1}.$$

将此式及(4.7)代入(3.6), 则得

$$\begin{aligned} N_n(x_0 + th) &= f_0 + t f_0 + \frac{t(t-1)}{2!}^2 f_0 + \dots \\ &\quad + \frac{t(t-1)\dots(t-n+1)}{n!}^n f_0, \end{aligned} \quad (4.10)$$

称为牛顿前插公式, 其余项由(2.14)得

$$R_n(x) = \frac{t(t-1)\dots(t-n)}{(n+1)!} h^{n+1} f^{(n+1)}(\), \quad (x_0, x_n). \quad (4.11)$$

如果要求函数表示 x_n 附近的函数值 $f(x)$, 此时应用牛顿插值公式(3.6), 插值点应按 x_n, x_{n-1}, \dots, x_0 的次序排列, 有

$$\begin{aligned} N_n(x) &= f(x_n) + f[x_n, x_{n-1}](x - x_n) \\ &\quad + f[x_n, x_{n-1}, x_{n-2}](x - x_n)(x - x_{n-1}) + \dots \\ &\quad + f[x_n, x_{n-1}, \dots, x_0](x - x_n)\dots(x - x_1). \end{aligned}$$

作变换 $x = x_n + th$ ($-1 \leq t \leq 0$), 并利用公式(4.8), 代入上式得

$$\begin{aligned} N_n(x_n + th) &= f_n + t f_n + \frac{t(t+1)}{2!}^2 f_n + \dots \\ &\quad + \frac{t(t+1)\dots(t+n-1)}{n!}^n f_n, \end{aligned} \quad (4.12)$$

称其为牛顿后插公式, 其余项

$$\begin{aligned} R_n(x) &= f(x) - N_n(x_n + th) \\ &= \frac{t(t+1)\dots(t+n)h^{n+1}}{(n+1)!} f^{(n+1)}(\), \end{aligned} \quad (4.13)$$

其中 (x_0, x_n) .

通常求开头部分插值点附近函数值时使用牛顿前插公式, 求插值节点末尾附近函数值时使用牛顿后插公式. 如果用相同节点进行插值, 则向前向后两种公式只是形式上差别, 其计算结果是相同的.

例 3 给出 $f(x) = \cos x$ 在 $x_k = kh, k = 0, 1, \dots, 6, h = 0.1$ 处的函数值, 试用 4 次等距节点插值公式计算 $f(0.048)$ 及 $f(0.566)$ 的近似值并估计误差.

解 先构造差分表. 用牛顿向前插值公式 (4.10) 计算 $f(0.048)$ 的近似值, 取 $x = 0.048, h = 0.1, t = \frac{x-0}{h} = 0.48$, 用表 2-4 上半部差分, 得

表 2-4

$f(x_k)$	$f(f)$	${}^2f({}^2f)$	${}^3f({}^3f)$	${}^4f({}^4f)$	${}^5f({}^5f)$
<u>1.00000</u>					
	<u>-0.00500</u>				
0.99500		<u>-0.00993</u>			
	-0.01493		<u>0.00013</u>		
0.98007		-0.00980		<u>0.00012</u>	
	-0.02473		0.00025		-0.00002
0.95534		-0.00955		0.00010	
	-0.03428		0.00035		-0.00001
0.92106		-0.00920		<u>0.00009</u>	

续表

$f(x_k)$	$f(f)$	${}^2 f({}^2 f)$	${}^3 f({}^3 f)$	${}^4 f({}^4 f)$	${}^5 f({}^5 f)$
0.87758	- 0.04348	<u>- 0.00876</u>	<u>0.00044</u>		
<u>0.82534</u>	<u>- 0.05224</u>				

$$\begin{aligned}
 N_4(0.048) &= 1.00000 + 0.48 \times (-0.00500) \\
 &\quad + \frac{(0.48)(0.48 - 1)}{2} (-0.00993) \\
 &\quad + \frac{1}{3!}(0.48)(0.48 - 1)(0.48 - 2)(0.00013) \\
 &\quad + \frac{1}{4!}(0.48)(0.48 - 1)(0.48 - 2)(0.48 - 3)(0.00012) \\
 &= 0.99885 \cos 0.048,
 \end{aligned}$$

误差估计由(4.11)可得

$$\begin{aligned}
 |R_4(0.048)| &\leq \frac{M_5}{5!} |t(t-1)(t-2)(t-3)(t-4)| / h^5 \\
 &\leq 1.5845 \times 10^{-7},
 \end{aligned}$$

其中 $M_5 = |\sin 0.6| = 0.565$.

计算 $f(0.566)$. 可用牛顿向后插值公式(4.12), $x = 0.566$,

$x_6 = 0.6$, $t = \frac{x - x_6}{h} = -0.34$, 用差分表 2-4 中下半部差分, 得

$$\begin{aligned}
 N_4(0.566) &= 0.82534 - 0.34 - 0.05224 + (0.66) \frac{-0.00876}{2} \\
 &\quad + (1.66) \frac{0.00044}{6} + 2.66 \times \frac{0.00009}{24} \\
 &= 0.84405.
 \end{aligned}$$

于是 $\cos 0.566 = 0.84405$, 误差估计由(4.13)得

$$\begin{aligned} & |R_4(0.566)| = \frac{M_5}{5!} |t(t+1)(t+2)(t+3)(t+4)| / h^5 \\ & \quad 1.7064 \times 10^{-7}, \end{aligned}$$

其中 $M_5 = 0.565$.

2.5 埃尔米特插值

不少实际的插值问题不但要求在节点上函数值相等, 而且还要求对应的导数值也相等, 甚至要求高阶导数也相等, 满足这种要求的插值多项式就是埃尔米特(Hermite)插值多项式. 下面只讨论函数值与导数值个数相等的情况. 设在节点 $a = x_0 < x_1 < \dots < x_n = b$ 上, $y_j = f(x_j)$, $m_j = f'(x_j)$ ($j = 0, 1, \dots, n$), 要求插值多项式 $H(x)$, 满足条件

$$H(x_j) = y_j, \quad H'(x_j) = m_j \quad (j = 0, 1, \dots, n). \quad (5.1)$$

这里给出了 $2n+2$ 个条件, 可唯一确定一个次数不超过 $2n+1$ 的多项式 $H_{2n+1}(x) = H(x)$, 其形式为

$$H_{2n+1}(x) = a_0 + a_1 x + \dots + a_{2n+1} x^{2n+1},$$

如根据条件(5.1)来确定 $2n+2$ 个系数 $a_0, a_1, \dots, a_{2n+1}$, 显然非常复杂, 因此, 我们仍采用求拉格朗日插值多项式的基函数方法. 先求插值基函数 $\varphi_j(x)$ 及 $\psi_j(x)$ ($j = 0, 1, \dots, n$), 共有 $2n+2$ 个, 每一个基函数都是 $2n+1$ 次多项式, 且满足条件

$$\varphi_j(x_k) = \psi_{jk} = \begin{cases} 0, & j \neq k, \\ 1, & j = k, \end{cases} \quad \varphi_j(x_k) = 0; \quad (5.2)$$

$$\varphi_j(x_k) = 0, \quad \psi_j(x_k) = \psi_{jk} (j, k = 0, 1, \dots, n).$$

于是满足条件(5.1)的插值多项式 $H(x) = H_{2n+1}(x)$ 可写成用插值基函数表示的形式

$$H_{2n+1}(x) = \sum_{j=0}^n [y_j \varphi_j(x) + m_j \psi_j(x)]. \quad (5.3)$$

由条件(5.2), 显然有 $H_{2n+1}(x_k) = y_k$, $H_{2n+1}(x_k) = m_k$, ($k=0, 1, \dots, n$). 下面的问题就是求满足条件(5.2)的基函数 ${}_j(x)$ 及 ${}_j(x)$. 为此, 可利用拉格朗日插值基函数 $l_j(x)$. 令

$${}_j(x) = (ax + b) \hat{l}_j(x),$$

其中 $l_j(x)$ 是(2.8)所表示的基函数. 由条件(5.2)有

$${}_j(x_j) = (ax_j + b) \hat{l}_j(x_j) = 1,$$

$${}_j(x_j) = l_j(x_j) [al_j(x_j) + 2(ax_j + b)l_j(x_j)] = 0,$$

$$ax_j + b = 1;$$

整理得

$$a + 2l_j(x_j) = 0.$$

解出

$$a = -2l_j(x_j), \quad b = 1 + 2x_j l_j(x_j).$$

由于

$$l_j(x) = \frac{(x - x_0) \dots (x - x_{j-1})(x - x_{j+1}) \dots (x - x_n)}{(x_j - x_0) \dots (x_j - x_{j-1})(x_j - x_{j+1}) \dots (x_j - x_n)},$$

利用两端取对数再求导, 得

$$l_j(x_j) = \prod_{k=0, k \neq j}^n \frac{1}{x_j - x_k},$$

于是

$${}_j(x) = 1 - 2(x - x_j) \prod_{k=0, k \neq j}^n \frac{1}{x_j - x_k} \hat{l}_j(x). \quad (5.4)$$

同理, 可得

$${}_j(x) = (x - x_j) \hat{l}_j(x). \quad (5.5)$$

还可证明满足条件(5.1)的插值多项式是唯一的. 用反证法, 假设 $H_{2n+1}(x)$ 及 $\tilde{H}_{2n+1}(x)$ 均满足条件(5.1), 于是

$$(x) = H_{2n+1}(x) - \tilde{H}_{2n+1}(x)$$

在每个节点 x_k 上均有二重根, 即 (x) 有 $2n+2$ 重根. 但 (x) 是不高于 $2n+1$ 次的多项式, 故 $(x) = 0$. 唯一性得证.

仿照拉格朗日插值余项的证明方法, 若 $f(x)$ 在 (a, b) 内的 $2n+2$ 阶导数存在, 则其插值余项

$$R(x) = f(x) - H_{2n+1}(x) = \frac{f^{(2n+2)}(\cdot)}{(2n+2)!} {}_{n+1}^2(x), \quad (5.6)$$

其中 (a, b) 且与 x 有关. 具体证明请读者自行完成.

作为带导数插值多项式(5.3)的重要特例是 $n=1$ 的情形. 这时可取节点为 x_k 及 x_{k+1} , 插值多项式为 $H_3(x)$, 满足条件

$$\begin{aligned} H_3(x_k) &= y_k, & H_3(x_{k+1}) &= y_{k+1}; \\ H_3(x_k) &= m_k, & H_3(x_{k+1}) &= m_{k+1}. \end{aligned} \quad (5.7)$$

相应的插值基函数为 ${}_k(x)$, ${}_{k+1}(x)$, ${}_k(x)$, ${}_{k+1}(x)$, 它们满足条件

$$\begin{aligned} {}_k(x_k) &= 1, & {}_k(x_{k+1}) &= 0, \\ {}_k(x_k) &= {}_k(x_{k+1}) = 0, \\ {}_{k+1}(x_k) &= 0, & {}_{k+1}(x_{k+1}) &= 1, \\ {}_{k+1}(x_k) &= {}_{k+1}(x_{k+1}) = 0; \\ {}_k(x_k) &= {}_k(x_{k+1}) = 0, \\ {}_k(x_k) &= 1, & {}_k(x_{k+1}) &= 0, \\ {}_{k+1}(x_k) &= {}_{k+1}(x_{k+1}) = 0, \\ {}_{k+1}(x_k) &= 0, & {}_{k+1}(x_{k+1}) &= 1. \end{aligned}$$

根据(5.4)及(5.5)的一般表达式, 可得到

$${}_k(x) = 1 + 2 \frac{x - x_k}{x_{k+1} - x_k} \frac{\frac{x - x_{k+1}}{x_k - x_{k+1}}^2}, \quad (5.8)$$

$${}_{k+1}(x) = 1 + 2 \frac{x - x_{k+1}}{x_k - x_{k+1}} \frac{\frac{x - x_k}{x_{k+1} - x_k}^2}.$$

$${}_k(x) = (x - x_k) \frac{\frac{x - x_{k+1}}{x_k - x_{k+1}}^2}, \quad (5.9)$$

$${}_{k+1}(x) = (x - x_{k+1}) \frac{\frac{x - x_k}{x_{k+1} - x_k}^2}.$$

于是满足条件(5.7)的插值多项式是

$$\begin{aligned} H_3(x) &= y_{k-k}(x) + y_{k+1-k+1}(x) + m_{k-k}(x) \\ &\quad + m_{k+1-k+1}(x), \end{aligned} \tag{5.10}$$

其余项 $R_3(x) = f(x) - H_3(x)$, 由(5.6)得

$$R_3(x) = \frac{1}{4!} f^{(4)}(\cdot)(x - x_k)^2(x - x_{k+1})^2, \quad (x_k, x_{k+1}).$$

例4 求满足 $P(x_j) = f(x_j)$ ($j = 0, 1, 2$) 及 $P(x_1) = f(x_1)$ 的插值多项式及其余项表达式.

由给定条件, 可确定次数不超过3的插值多项式. 由于此多项式通过点 $(x_0, f(x_0))$, $(x_1, f(x_1))$ 及 $(x_2, f(x_2))$, 故其形式为

$$\begin{aligned} P(x) &= f(x_0) + f[x_0, x_1](x - x_0) \\ &\quad + f[x_0, x_1, x_2](x - x_0)(x - x_1) \\ &\quad + A(x - x_0)(x - x_1)(x - x_2), \end{aligned}$$

其中 A 为待定常数, 可由条件 $P(x_1) = f(x_1)$ 确定, 通过计算可得

$$A = \frac{f(x_1) - f[x_0, x_1] - (x_1 - x_0)f[x_0, x_1, x_2]}{(x_1 - x_0)(x_1 - x_2)}.$$

为了求出余项 $R(x) = f(x) - P(x)$ 的表达式, 可设

$$R(x) = f(x) - P(x) = k(x)(x - x_0)(x - x_1)^2(x - x_2),$$

其中 $k(x)$ 为待定函数. 构造

$$(t) = f(t) - P(t) - k(x)(t - x_0)(t - x_1)^2(t - x_2).$$

显然 $(x_j) = 0$ ($j = 0, 1, 2$). 且 $(x_1) = 0$, $(x) = 0$, 故 (t) 在 (a, b) 内有5个零点(二重根算两个). 反复应用罗尔定理, 得 ${}^{(4)}(t)$ 在 (a, b) 内至少有一个零点, 故

$${}^{(4)}(\cdot) = f^{(4)}(\cdot) - 4!k(x) = 0,$$

于是

$$k(x) = \frac{1}{4!} f^{(4)}(\cdot),$$

余项表达式为

$$R(x) = \frac{1}{4!} f^{(4)}(\xi)(x - x_0)(x - x_1)^2(x - x_2), \quad (5.11)$$

式中 ξ 位于 x_0, x_1, x_2 和 x 所界定的范围内。

2.6 分段低次插值

2.6.1 高次插值的病态性质

上面我们根据区间 $[a, b]$ 上给出的节点做插值多项式 $L_n(x)$ 近似 $f(x)$, 一般总认为 $L_n(x)$ 的次数 n 越高逼近 $f(x)$ 的精度越好, 但实际上并非如此。这是因为对任意的插值节点, 当 n 时, $L_n(x)$ 不一定收敛到 $f(x)$ 。20 世纪初龙格(Runge)就给出了一个等距节点插值多项式 $L_n(x)$ 不收敛到 $f(x)$ 的例子。他给出的函数为 $f(x) = V(1 + x^2)$, 它在 $[-5, 5]$ 上各阶导数均存在。在 $[-5, 5]$ 上取 $n+1$ 个等距节点 $x_k = -5 + 10 \frac{k}{n}$ ($k=0, 1, \dots, n$) 所构造的拉格朗日插值多项式为

$$L_n(x) = \sum_{j=0}^n \frac{1}{1+x_j^2} \frac{\frac{n+1}{n+1}(x)}{(x-x_j)\frac{n+1}{n+1}(x_j)}.$$

令 $x_{n+1/2} = \frac{1}{2}(x_{n+1} + x_n)$, 则 $x_{n+1/2} = 5 - \frac{5}{n}$, 表 2-5 列出了 $n=2, 4, \dots, 20$ 的 $L_n(x_{n+1/2})$ 的计算结果及在 $x_{n+1/2}$ 上的误差 $R(x_{n+1/2})$ 。可以看出, 随 n 的增加, $R(x_{n+1/2})$ 的绝对值几乎成倍地增加。这说明当 n 时 L_n 在 $[-5, 5]$ 上不收敛。Runge 证明了, 存在一个常数 $c=3.63$, 使得当 $|x| > c$ 时, $\lim_n L_n(x) = f(x)$, 而当 $|x| > c$ 时 $\{L_n(x)\}$ 发散。

下面取 $n=10$, 根据计算画出 $y=L_{10}(x)$ 及 $y=V(1+x^2)$ 在 $[-5, 5]$ 上的图形, 见图 2-5。

表 2-5

n	$f(x_{n-1/2})$	$L_n(x_{n-1/2})$	$R(x_{n-1/2})$
2	0.137931	0.759615	-0.621684
4	0.066390	-0.356826	0.423216
6	0.054463	0.607879	-0.553416
8	0.049651	-0.831017	0.880668
10	0.047059	1.578721	-1.531662
12	0.045440	-2.755000	2.800440
14	0.044334	5.332743	-5.288409
16	0.043530	-10.173867	10.217397
18	0.042920	20.123671	-20.080751
20	0.042440	-39.952449	39.994889

从图上看到, 在 $x = \pm 5$ 附近 $L_{10}(x)$ 与 $f(x) = 1/(1+x^2)$ 偏离很远, 例如 $L_{10}(4.8) = 1.80438$, $f(4.8) = 0.04160$. 这说明用高次插值多项式 $L_n(x)$ 近似 $f(x)$ 效果并不好, 因而通常不用高次插值, 而用分段低次插值. 从本例看到, 如果我们把 $y = 1/(1+x^2)$ 在节点 $x = 0, \pm 1, \pm 2, \pm 3, \pm 4, \pm 5$ 处用折线连起来显然比 $L_{10}(x)$ 逼近 $f(x)$ 好得多. 这正是我们下面要讨论的分段低次插值的出发点.

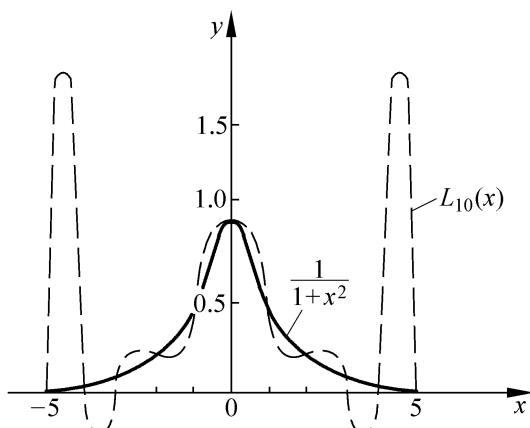


图 2-5

2.6.2 分段线性插值

所谓分段线性插值就是通过插值点用折线段连接起来逼近 $f(x)$ 。设已知节点 $a = x_0 < x_1 < \dots < x_n = b$ 上的函数值 f_0, f_1, \dots, f_n , 记 $h_k = x_{k+1} - x_k$, $h = \max_k h_k$, 求一折线函数 $I_h(x)$ 满足:

1° 记 $I_h(x) \in C[a, b]$,

2° $I_h(x_k) = f_k$ ($k = 0, 1, \dots, n$),

3° $I_h(x)$ 在每个小区间 $[x_k, x_{k+1}]$ 上是线性函数。

则称 $I_h(x)$ 为分段线性插值函数。

由定义可知 $I_h(x)$ 在每个小区间 $[x_k, x_{k+1}]$ 上可表示为

$$I_h(x) = \frac{x - x_{k+1}}{x_k - x_{k+1}} f_k + \frac{x - x_k}{x_{k+1} - x_k} f_{k+1} \quad (x_k \leq x \leq x_{k+1}) . \quad (6.1)$$

若用插值基函数表示, 则在整个区间 $[a, b]$ 上 $I_h(x)$ 为

$$I_h(x) = \sum_{j=0}^n f_j l_j(x), \quad (6.2)$$

其中基函数 $l_j(x)$ 满足条件 $l_j(x_k) = \delta_{jk}$ ($j, k = 0, 1, \dots, n$), 其形式是

$$\frac{x - x_{j-1}}{x_j - x_{j-1}}, \quad x_{j-1} \leq x \leq x_j \quad (j = 0 \text{ 略去});$$

$$l_j(x) = \frac{x - x_{j+1}}{x_j - x_{j+1}}, \quad x_j \leq x \leq x_{j+1} \quad (j = n \text{ 略去}); \quad (6.3)$$

$$0, \quad x \in [a, b], x \notin [x_{j-1}, x_{j+1}] .$$

分段线性插值的误差估计可利用插值余项(2.17)得到

$$\max_{x_k \leq x \leq x_{k+1}} |f(x) - I_h(x)| \leq \frac{M_2}{2} \max_{x_k \leq x \leq x_{k+1}} |(x - x_k)(x - x_{k+1})|$$

或

$$\max_a \max_b |f(x) - I_h(x)| \leq \frac{M_2}{8} h^2, \quad (6.4)$$

其中 $M_2 = \max_{a \leq x \leq b} |f(x)|$.

分段线性插值基函数 $l_j(x)$ 只在 x_j 附近不为零, 在其他地方均为零, 这种性质称为局部非零性质. 当 $x \in [x_k, x_{k+1}]$ 时

$$1 = \sum_{j=0}^n l_j(x) = l_k(x) + l_{k+1}(x),$$

故

$$f(x) = [l_k(x) + l_{k+1}(x)] f(x).$$

另一方面, 这时

$$I_h(x) = f_k l_k(x) + f_{k+1} l_{k+1}(x).$$

现在证明 $\lim_{h \rightarrow 0} I_h(x) = f(x)$. 考虑

$$\begin{aligned} & |f(x) - I_h(x)| = |l_k(x)| |f(x) - f_k| \\ & + |l_{k+1}(x)| |f(x) - f_{k+1}| \\ & [l_k(x) + l_{k+1}(x)] (h_k) = (h). \end{aligned}$$

这里 (h) 是函数 $f(x)$ 在区间 $[a, b]$ 上的连续模, 即对任意两点 $x, x \in [a, b]$, 只要 $|x - x| \leq h$, 就有

$$|f(x) - f(x)| \leq (h)$$

称 (h) 为 $f(x)$ 在 $[a, b]$ 上的连续模, 当 $f(x) \in C[a, b]$ 时, 就有 $\lim_{h \rightarrow 0} (h) = 0$.

由前式可知, 当 $x \in [a, b]$ 时有

$$\max_{a \leq x \leq b} |f(x) - I_h(x)| \leq (h).$$

因此, 只要 $f(x) \in C[a, b]$, 就有

$$\lim_{h \rightarrow 0} I_h(x) = f(x)$$

在 $[a, b]$ 上一致成立, 故 $I_h(x)$ 在 $[a, b]$ 上一致收敛到 $f(x)$.

2.6.3 分段三次埃尔米特插值

分段线性插值函数 $I_h(x)$ 的导数是间断的, 若在节点 x_k ($k = 0, 1, \dots, n$) 上除已知函数值 f_k 外还给出导数值 $f'_k = m_k$ ($k = 0, 1, \dots, n$), 这样就可构造一个导数连续的分段插值函数 $I_h(x)$, 它

满 条件:

1. $I_h(x) \in C[a, b]$ ($C[a, b]$ 代表区间 $[a, b]$ 上一阶导数连续的函数集合),

2. $I_h(x_k) = f_k$, $I_h(x_{k+1}) = f_{k+1}$ ($k = 0, 1, \dots, n$),

3. $I_h(x)$ 在每个小区间 $[x_k, x_{k+1}]$ 上是三次多项式.

根据两点三次插值多项式(5.10). 可知, $I_h(x)$ 在区间 $[x_k, x_{k+1}]$ 上的表达式为

$$\begin{aligned} I_h(x) &= \frac{x - x_{k+1}}{x_k - x_{k+1}}^2 \cdot 1 + 2 \frac{x - x_k}{x_{k+1} - x_k} f_k + \frac{x - x_k}{x_{k+1} - x_k}^2 \\ &\quad \cdot 1 + 2 \frac{x - x_{k+1}}{x_k - x_{k+1}} f_{k+1} + \frac{x - x_{k+1}}{x_k - x_{k+1}}^2 (x - x_k) f_k \\ &\quad + \frac{x - x_k}{x_{k+1} - x_k}^2 (x - x_{k+1}) f_{k+1}. \end{aligned} \quad (6.5)$$

若在整个区间 $[a, b]$ 上定义一组分段三次插值基函数 $\varphi_j(x)$ 及 $\psi_j(x)$ ($j = 0, 1, \dots, n$), 则 $I_h(x)$ 可表示为

$$I_h(x) = \sum_{j=0}^n [f_j \varphi_j(x) + f_{j+1} \psi_j(x)], \quad (6.6)$$

其中 $\varphi_j(x)$, $\psi_j(x)$ 分别表示为

$$\frac{x - x_{j-1}}{x_j - x_{j-1}}^2 \cdot 1 + 2 \frac{x - x_j}{x_{j-1} - x_j}, \quad \begin{matrix} x_{j-1} & x & x_j \\ (j=0 \text{ 略去}) \end{matrix}$$

$$\varphi_j(x) = \frac{x - x_{j+1}}{x_j - x_{j+1}}^2 \cdot 1 + 2 \frac{x - x_j}{x_{j+1} - x_j}, \quad \begin{matrix} x_j & x & x_{j+1} \\ (j=n \text{ 略去}) \end{matrix} \quad (6.7)$$

$$0, \quad \text{其他};$$

$$\frac{x - x_{j-1}}{x_j - x_{j-1}}^2 (x - x_j), \quad \begin{matrix} x_{j-1} & x & x_j \\ (j=0 \text{ 略去}) \end{matrix}$$

$$\varphi_j(x) = \frac{x - x_{j+1}}{x_j - x_{j+1}}^2 (x - x_j), \quad \begin{matrix} x_j & x & x_{j+1} \\ (j=n \text{ 略去}) \end{matrix}$$

$$0, \quad \text{其他}.$$

$$(6.8)$$

由于 $j(x), j(x)$ 的局部非零性质, 当 $x \in [x_k, x_{k+1}]$ 时, 只有 $k(x), k_{+1}(x), k(x), k_{+1}(x)$ 不为零, 于是(6.6)的 $I_h(x)$ 可表示为

$$I_h(x) = f_{k-k}(x) + f_{k+1-k+1}(x) + f_{k-k}(x) + f_{k+1-k+1}(x) \\ (x_k \leq x \leq x_{k+1}). \quad (6.9)$$

为了研究 $I_h(x)$ 的收敛性, 由(6.7)及(6.8)直接得估计式

$$0 \leq j(x) \leq 1, \quad (6.10)$$

$$|k(x)| \leq \frac{4}{27} h_k, \quad |k_{+1}(x)| \leq \frac{4}{27} h_k. \quad (6.11)$$

此外, 当 $f(x)$ 是分段三次多项式时, $f(x)$ 的插值多项式 $I_h(x)$ 就是它本身. 例如, 当 $f(x) = 1$ 时就有

$$\sum_{j=0}^n j(x) = 1.$$

当 $x \in [x_k, x_{k+1}]$ 时, 就得

$$k(x) + k_{+1}(x) = 1. \quad (6.12)$$

由(6.9)~(6.12), 当 $x \in [x_k, x_{k+1}]$ 时还可得

$$|f(x) - I_h(x)| \leq |k(x)| |f(x) - f_k| + |k_{+1}(x)| |f(x) - f_{k+1}| \\ + (4/27) h_k [|f_k| + |f_{k+1}|] \\ |k(x)| |f(\cdot)| h_k + |k_{+1}(x)| |f(\cdot)| h_{k+1} \\ + \frac{4}{27} [|f_k| + |f_{k+1}|] h_k,$$

这里, (x_k, x_{k+1}) 且依赖于 x . 因此对 $x \in [a, b]$ 成立

$$\max_{a \leq x \leq b} |f(x) - I_h(x)| \leq \frac{35}{27} h \max_{a \leq x \leq b} |f(x)|. \quad (6.13)$$

这表明用 $I_h(x)$ 逼近 $f(x)$ 时, 它的界只依赖 h , 而与 x 无关. 因此, 当 $x \in [a, b]$ 时

$$\lim_{h \rightarrow 0} I_h(x) = f(x)$$

一致成立. 从而得到:

定理3 设 $f \in C[a, b]$, 则当 $h \rightarrow 0$ 时, $I_h(x)$ 在 $[a, b]$ 上一致

收敛于 $f(x)$.

2.7 三次样条插值

上面讨论的分段低次插值函数都有一致收敛性, 但光滑性较差, 对于像高速飞机的机翼形线, 船体放样等型值线往往要求有二阶光滑度, 即有二阶连续导数. 早期工程师制图时, 把富有弹性的细长木条(所谓样条)用压铁固定在样点上, 在其他地方让它自由弯曲, 然后画下长条的曲线, 称为样条曲线. 样条曲线实际上是由分段三次曲线并接而成, 在连接点即样点上要求二阶导数连续, 从数学上加以概括就得到数学样条这一概念. 下面我们讨论最常用的三次样条函数.

2.7.1 三次样条函数

定义 4 若函数 $S(x) \in C^2[a, b]$, 且在每个小区间 $[x_j, x_{j+1}]$ 上是三次多项式, 其中 $a = x_0 < x_1 < \dots < x_n = b$ 是给定节点, 则称 $S(x)$ 是节点 x_0, x_1, \dots, x_n 上的三次样条函数. 若在节点 x_j 上给定函数值 $y_j = f(x_j)$ ($j = 0, 1, \dots, n$), 并成立

$$S(x_j) = y_j \quad (j = 0, 1, \dots, n), \quad (7.1)$$

则称 $S(x)$ 为三次样条插值函数.

从定义知要求出 $S(x)$, 在每个小区间 $[x_j, x_{j+1}]$ 上要确定 4 个待定系数, 共有 n 个小区间, 故应确定 $4n$ 个参数. 根据 $S(x)$ 在 $[a, b]$ 上二阶导数连续, 在节点 x_j ($j = 1, 2, \dots, n - 1$) 处应满足连续性条件

$$\begin{aligned} S(x_j - 0) &= S(x_j + 0), \quad S'(x_j - 0) = S'(x_j + 0), \\ S''(x_j - 0) &= S''(x_j + 0). \end{aligned} \quad (7.2)$$

共有 $3n - 3$ 个条件, 再加上 $S(x)$ 满足插值条件 (7.1), 共有 $4n - 2$ 个条件, 因此还需要 2 个条件才能确定 $S(x)$. 通常可在区间 $[a, b]$

端点 $a = x_0$, $b = x_n$ 上各加一个条件(称为边界条件), 可根据实际问题的要求给定. 常见的有以下3种:

1° 已知两端的一阶导数值, 即

$$S'(x_0) = f_0, \quad S'(x_n) = f_n. \quad (7.3)$$

2° 两端的二阶导数已知, 即

$$S''(x_0) = f_0, \quad S''(x_n) = f_n, \quad (7.4)$$

其特殊情况为

$$S(x_0) = S(x_n) = 0. \quad (7.4)$$

(7.4) 称为自然边界条件.

3° 当 $f(x)$ 是以 $x_n - x_0$ 为周期的周期函数时, 则要求 $S(x)$ 也是周期函数. 这时边界条件应满足

$$S(x_0 + 0) = S(x_n - 0), \quad S'(x_0 + 0) = S'(x_n - 0),$$

$$S''(x_0 + 0) = S''(x_n - 0), \quad (7.5)$$

而此时(7.1)中 $y_0 = y_n$. 这样确定的样条函数 $S(x)$ 称为周期样条函数.

2.7.2 样条插值函数的建立

构造满足插值条件(7.1)及相应边界条件的三次样条插值函数 $S(x)$ 的表达式可以有多种方法. 例如, 可以直接利用分段三次埃尔米特插值(6.6), 只要假定 $S(x_j) = m_j$ ($j = 0, 1, \dots, n$), 再由(7.1)可得

$$S(x) = \sum_{j=0}^n [y_{j-j}(x) + m_{j-j}(x)], \quad (7.6)$$

其中 $\psi_j(x)$, $\varphi_j(x)$ 是由(6.7), (6.8)表示的插值基函数, 利用条件(7.2)及相应边界条件(7.3)~(7.5)则可得到关于 m_j ($j = 0, 1, \dots, n$) 的三对角方程组, 求出 m_j 则得到所求的三次样条函数 $S(x)$.

下面我们利用 $S(x)$ 的二阶导数值 $S''(x_j) = M_j$ ($j = 0, 1, \dots,$

n)表达 $S(x)$, 由于 $S(x)$ 在区间 $[x_j, x_{j+1}]$ 上是三次多项式, 故 $S(x)$ 在 $[x_j, x_{j+1}]$ 上是线性函数, 可表示为

$$S(x) = M_j \frac{x_{j+1} - x}{h_j} + M_{j+1} \frac{x - x_j}{h_j}. \quad (7.7)$$

对 $S(x)$ 积分两次并利用 $S(x_j) = y_j$ 及 $S(x_{j+1}) = y_{j+1}$, 可定出积分常数, 于是得三次样条表达式

$$\begin{aligned} S(x) &= M_j \frac{(x_{j+1} - x)^3}{6 h_j} + M_{j+1} \frac{(x - x_j)^3}{6 h_j} \\ &+ y_j - \frac{M_j h_j^2}{6} \frac{x_{j+1} - x}{h_j} + y_{j+1} - \frac{M_{j+1} h_j^2}{6} \frac{x - x_j}{h_j} \\ &\quad (j = 0, 1, \dots, n - 1). \end{aligned} \quad (7.8)$$

这里 $M_j, j = 0, 1, \dots, n$, 是未知的. 为了确定 $M_j, j = 0, 1, \dots, n$, 对 $S(x)$ 求导得

$$\begin{aligned} S'(x) &= -M_j \frac{(x_{j+1} - x)^2}{2 h_j} + M_{j+1} \frac{(x - x_j)^2}{2 h_j} \\ &+ \frac{y_{j+1} - y_j}{h_j} - \frac{M_{j+1} - M_j}{6} h_j; \end{aligned} \quad (7.9)$$

由此可求得

$$S(x_j + 0) = -\frac{h_j}{3} M_j - \frac{h_j}{6} M_{j+1} + \frac{y_{j+1} - y_j}{h_j}.$$

类似地可求出 $S(x)$ 在区间 $[x_{j-1}, x_j]$ 上的表达式, 从而得

$$S(x_j - 0) = \frac{h_{j-1}}{6} M_{j-1} + \frac{h_{j-1}}{3} M_j + \frac{y_j - y_{j-1}}{h_{j-1}},$$

利用 $S(x_j + 0) = S(x_j - 0)$ 可得

$$\mu_j M_{j-1} + 2M_j + \mu_j M_{j+1} = d_j \quad (j = 1, 2, \dots, n - 1), \quad (7.10)$$

其中

$$\mu_j = \frac{h_{j-1}}{h_{j-1} + h_j}, \quad \mu_j = \frac{h_j}{h_{j-1} + h_j}, \quad j = 0, 1, \dots, n,$$

$$d_j = 6 \frac{f[x_j, x_{j+1}] - f[x_{j-1}, x_j]}{h_{j-1} + h_j} = 6 f[x_{j-1}, x_j, x_{j+1}] . \quad (7.11)$$

对第一种边界条件(7.3), 可导出两个方程

$$\begin{aligned} 2M_0 + M_1 &= \frac{6}{h_0} (f[x_0, x_1] - f_0), \\ M_{n-1} + 2M_n &= \frac{6}{h_{n-1}} (f_n - f[x_{n-1}, x_n]). \end{aligned} \quad (7.12)$$

如果令 $\mu_0 = 1$, $d_0 = \frac{6}{h_0} (f[x_0, x_1] - f_0)$, $\mu_n = 1$, $d_n = \frac{6}{h_{n-1}} (f_n - f[x_{n-1}, x_n])$, 那么(7.10)及(7.12)可写成矩阵形式

$$\begin{matrix} 2 & & 0 & & & M_0 & d_0 \\ \mu & 2 & & 1 & & M_1 & d_1 \\ & W & & W & & \dots & \dots \\ & \mu_{n-1} & & 2 & & M_{n-1} & d_{n-1} \\ & \mu_n & & 2 & & M_n & d_n \end{matrix} = \dots \quad (7.13)$$

对第二种边界条件(7.4), 直接得端点方程

$$M_0 = f_0, \quad M_n = f_n. \quad (7.14)$$

如果令 $\mu_0 = \mu_n = 0$, $d_0 = 2f_0$, $d_n = 2f_n$, 则(7.10)和(7.14)也可以写成(7.13)的形式.

对于第三种边界条件(7.5), 可得

$$M_0 = M_n, \quad \mu_n M_{n-1} + 2M_n = d_n, \quad (7.15)$$

其中 $\mu_n = \frac{h_0}{h_{n-1} + h_0}$, $\mu_n = 1 - \mu_n = \frac{h_{n-1}}{h_{n-1} + h_0}$, $d_n =$

$6 \frac{f[x_0, x_1] - f[x_{n-1}, x_n]}{h_0 + h_{n-1}}$, (7.10)和(7.15)可以写成矩阵形式

$$\begin{matrix} 2 & & 1 & & & \mu & M_1 & d_1 \\ \mu & 2 & & 2 & & M_2 & d_2 \\ & W & & W & & \dots & \dots \\ & \mu_{n-1} & & 2 & & M_{n-1} & d_{n-1} \\ & \mu_n & & 2 & & M_n & d_n \end{matrix} = \dots \quad (7.16)$$

(7.13)和(7.16)是关于 M_j ($j=0, 1, \dots, n$)的三对角方程组, M_j 在力学上解释为细梁在 x_j 截面处的弯矩, 称为 $S(x)$ 的矩, 方程组(7.13)和(7.16)称为三弯矩方程。(7.13)和(7.16)的系数矩阵中元素 μ_j 已完全确定。并且满足 $\mu_0 = 0, \mu_n = 0, \mu_{j-1} + \mu_j = 1$ 。因此系数矩阵为严格对角占优阵, 从而(7.13)和(7.16)有唯一解。求解方法可见第5章第4节追赶法, 将解得结果代入(7.8)的表达式即可。

例5 设 $f(x)$ 为定义在 $[27.7, 30]$ 上的函数, 在节点 x_i ($i=0, 1, 2, 3$) 上的值如下:

$$f(x_0) = f(27.7) = 4.1, f(x_1) = f(28) = 4.3,$$

$$f(x_2) = f(29) = 4.1, f(x_3) = f(30) = 3.0.$$

试求三次样条函数 $S(x)$, 使它满足边界条件 $S(27.7) = 3.0$, $S(30) = -4.0$ 。

解 先由(7.11)及(7.12)计算 $h_0 = 0.30, h_1 = h_2 = 1, \mu_0 = \frac{3}{13}, \mu_1 = \frac{1}{2}, \mu_2 = 1, \mu_3 = 1, M_0 = \frac{10}{13}, M_1 = \frac{1}{2}, d_0 = \frac{6}{h_0}(f[x_0, x_1] - f_0) = -46.666, d_1 = 6f[x_0, x_1, x_2] = -4.00002, d_2 = 6f[x_1, x_2, x_3] = -2.70000, d_3 = \frac{6}{h_2}(f_3 - f[x_2, x_3]) = -17.4$ 。

由此得矩阵形式的方程组(7.13)为

$$\begin{array}{ccccc} 2 & 1 & & & \\ \frac{3}{13} & 2 & \frac{10}{13} & M_0 & -46.6666 \\ & & & M_1 & -4.00002 \\ \frac{1}{2} & 2 & \frac{1}{2} & M_2 & -2.7000 \\ & & & M_3 & -17.4000 \\ 1 & 2 & & & \end{array} .$$

求解此方程组得到

$$M_0 = -23.531, M_1 = 0.395, M_2 = 0.830, M_3 = -9.115.$$

将 M_0, M_1, M_2, M_3 代入表达式(7.8)得到(曲线见图2-6)

$$\begin{aligned}
 & 13.07278(x - 28)^3 - 14.84322(x - 28) + 0.21944 \\
 & \quad \cdot (x - 27.7)^3 + 14.31358(x - 27.7), \quad x \in [27.7, 28], \\
 S(x) = & 0.06583(29 - x)^3 + 4.23417(29 - x) + 0.13833 \\
 & \quad \cdot (x - 28)^3 + 3.96167(x - 28), \quad x \in [28, 29], \\
 & 0.13833(30 - x)^3 + 3.96167(30 - x) - 1.51917 \\
 & \quad \cdot (x - 29)^3 + 4.51917(x - 29), \quad x \in [29, 30].
 \end{aligned}$$

通常求三次样条函数可根据上述例题的计算步骤直接编程上机计算, 或直接使用数学库中软件, 根据具体要求算出结果即可.

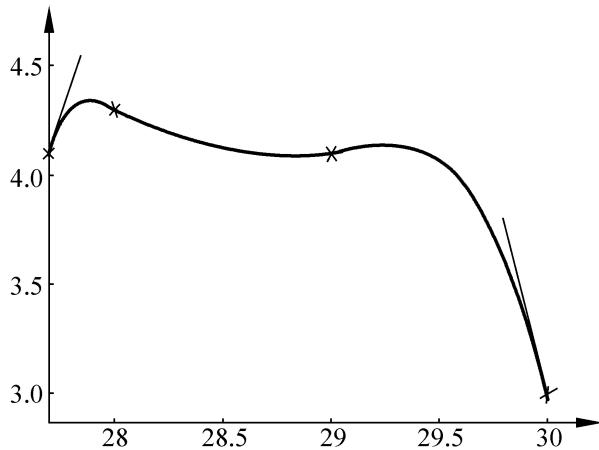


图 2-6

例 6 给定函数 $f(x) = \frac{1}{1+x^2}$, $-5 \leq x \leq 5$, 节点 $x_k = -5 + k$ ($k = 0, 1, \dots, 10$), 用三次样条插值求 $S_{10}(x)$.

取 $S_{10}(x_k) = f(x_k)$ ($k = 0, 1, \dots, 10$), $S_{10}(-5) = f(-5)$, $S_{10}(5) = f(5)$. 直接上机计算可求出 $S_{10}(x)$ 在表 2-6 所列各点的值. 从表中看到, 在所列各点 $S_{10}(x)$ 与 $f(x)$ 误差较小, 它可作为 $f(x)$ 在区间 $[-5, 5]$ 上的近似, 而用拉格朗日插值多项式 $L_{10}(x)$ 计算相应点上的值 $L_{10}(x)$ (也见表 2-6), 显然它与 $f(x)$ 相差很大, 在图 2-5 中已经看到它不能作为 $f(x)$ 的近似.

表 2-6

x	$\frac{1}{1+x^2}$	$S_{10}(x)$	$L_{10}(x)$	x	$\frac{1}{1+x^2}$	$S_{10}(x)$	$L_{10}(x)$
-5.0	0.03846	0.03846	0.03846	-2.3	0.15898	0.16115	0.24145
-4.8	0.04160	0.03758	1.80438	-2.0	0.20000	0.20000	0.20000
-4.5	0.04706	0.04248	1.57872	-1.8	0.23585	0.23154	0.18878
-4.3	0.05131	0.04842	0.88808	-1.5	0.30769	0.29744	0.23535
-4.0	0.05882	0.05882	0.05882	-1.3	0.37175	0.36133	0.31650
-3.8	0.06477	0.06556	-0.20130	-1.0	0.50000	0.50000	0.50000
-3.5	0.07547	0.07606	-0.22620	-0.8	0.60976	0.62420	0.64316
-3.3	0.08410	0.08426	-0.10832	-0.5	0.80000	0.82051	0.84340
-3.0	0.10000	0.10000	0.10000	-0.3	0.91743	0.92754	0.94090
-2.8	0.11312	0.11366	0.19837	0	1.00000	1.00000	1.00000
-2.5	0.13793	0.13971	0.25376				

2.7.3 误差界与收敛性

三次样条函数的收敛性与误差估计比较复杂, 这里不加证明地给出一个主要结果 .

定理4 设 $f(x) \in C^4[a, b]$, $S(x)$ 为满足第一种或第二种边界条件(7.3)或(7.4)的三次样条函数, 令 $h = \max_0^1 h_i$, $h_i = x_{i+1} - x_i$ ($i = 0, 1, \dots, n - 1$), 则有估计式

$$\max_a^b |f^{(k)}(x) - S^{(k)}(x)| \leq C_k \max_a^b |f^{(4)}(x)| h^{4-k}, \quad k = 0, 1, 2,$$
(7.17)

其中 $C_0 = \frac{5}{384}$, $C_1 = \frac{1}{24}$, $C_2 = \frac{3}{8}$.

这个定理不但给出了三次样条插值函数 $S(x)$ 的误差估计, 且当 $h \rightarrow 0$ 时, $S(x)$ 及其一阶导数 $S'(x)$ 和二阶导数 $S''(x)$ 均分别一

致收敛于 $f(x)$, $f'(x)$ 及 $f''(x)$.

评注

插值法是一个古老而实用的课题。它是函数逼近，数值微积分和微分方程数值解的基础。本章讨论了拉格朗日插值公式及牛顿插值公式，前者在理论上较为重要，后者在计算插值多次式及求函数近似值较为方便且节省计算量。等距节点插值是应用中最常见的，利用差分及牛顿前插与后插公式即可，还有利用中心差分得到的其他类型插值公式，因使用较少本章没有介绍，如有需要可参看文献[3]。

对充分光滑的被插函数可采用微分形式的误差估计给出误差限，其他情形可利用均差形式给出误差估计的近似值。但由于高次插值存在病态性质，一般实际计算很少使用高次插值，更多使用分段低次插值，特别是三次样条插值，由于它具有良好的收敛性和稳定性，又有二阶光滑度，因此在理论上和应用中均有重要意义，本章只对最常用的三弯矩方程做简单介绍，一般的样条函数理论和 B-样条等更详细内容读者如有需要可参看文献[4, 5, 10]等。

习题

1. 当 $x=1, -1, 2$ 时, $f(x)=0, -3, 4$, 求 $f(x)$ 的二次插值多项式。

2. 给出 $f(x) = \ln x$ 的数值表

x	0.4	0.5	0.6
$\ln x$	- 0.916291	- 0.693147	- 0.510826
x	0.7	0.8	
$\ln x$	- 0.356675	- 0.223144	

用线性插值及二次插值计算 $\ln 0.54$ 的近似值 .

3 . 给出 $\cos x, 0^\circ \leq x \leq 90^\circ$ 的函数表, 步长 $h = 1 = (1/60)^\circ$, 若函数表具有 5 位有效数字, 研究用线性插值求 $\cos x$ 近似值时的总误差界 .

4 . 设 x_j 为互异节点 ($j = 0, 1, \dots, n$), 求证:

$$\text{i)} \quad \sum_{j=0}^n x_j^k l_j(x) = x^k \quad (k = 0, 1, \dots, n);$$

$$\text{ii)} \quad \sum_{j=0}^n (x_j - x)^k l_j(x) = 0 \quad (k = 1, 2, \dots, n).$$

5 . 设 $f(x) \in C^2[a, b]$ 且 $f(a) = f(b) = 0$, 求证:

$$\max_{a \leq x \leq b} |f(x)| \leq \frac{1}{8}(b-a)^2 \max_{a \leq x \leq b} |f'(x)|.$$

6 . 在 $-4 \leq x \leq 4$ 上给出 $f(x) = e^x$ 的等距节点函数表, 若用二次插值求 e^x 的近似值, 要使截断误差不超过 10^{-6} , 问使用函数表的步长 h 应取多少?

7 . 若 $y_n = 2^n$, 求 4y_n 及 ${}^4\bar{y}_n$.

8 . 如果 $f(x)$ 是 m 次多项式, 记 $f(x) = f(x+h) - f(x)$, 证明 $f(x)$ 的 k 阶差分 ${}^k f(x) (0 \leq k \leq m)$ 是 $m-k$ 次多项式, 并且 ${}^{m+l} f(x) = 0$ (l 为正整数) .

9 . 证明 $(f_k g_k) = f_k g_k + g_{k+1} f_k$.

$$\text{10 . } \sum_{k=0}^{n-1} f_k g_k = f_n g_n - f_0 g_0 - \sum_{k=0}^{n-1} g_{k+1} f_k.$$

$$\text{11 . 证明 } \sum_{j=0}^{n-1} {}^2 y_j = y_n - y_0.$$

12 . 若 $f(x) = a_0 + a_1 x + \dots + a_{n-1} x^{n-1} + a_n x^n$ 有 n 个不同实根 x_1, x_2, \dots, x_n , 证明:

$$\sum_{j=1}^n \frac{x_j^k}{f(x_j)} = 0, 0 \leq k \leq n-2; \\ a_n^{-1}, \quad k = n-1.$$

13 . 证明 n 阶均差有下列性质:

i) 若 $F(x) = c f(x)$, 则 $F[x_0, x_1, \dots, x_n] = c f[x_0, x_1, \dots, x_n]$;

ii) 若 $F(x) = f(x) + g(x)$, 则

$$F[x_0, x_1, \dots, x_n] = f[x_0, x_1, \dots, x_n] + g[x_0, x_1, \dots, x_n].$$

14 . $f(x) = x^7 + x^4 + 3x + 1$, 求

$$f[2^0, 2^1, \dots, 2^7] \text{ 及 } f[2^0, 2^1, \dots, 2^8].$$

15. 证明两点三次埃尔米特插值余项是

$$R_3(x) = f^{(4)}(\cdot)(x - x_k)^2(x - x_{k+1})^2 / 4!, \quad (x_k, x_{k+1}),$$

并由此求出分段三次埃尔米特插值的误差限.

16. 求一个次数不高于4次的多项式 $P(x)$, 使它满足 $P(0) = P'(0) = 0$, $P(1) = P'(1) = 1$, $P(2) = 1$.

17. 设 $f(x) = 1/(1+x^2)$, 在 $-5 \leq x \leq 5$ 上取 $n=10$, 按等距节点求分段线性插值函数 $I_h(x)$, 计算各节点间中点处的 $I_h(x)$ 与 $f(x)$ 的值, 并估计误差.

18. 求 $f(x) = x^2$ 在 $[a, b]$ 上的分段线性插值函数 $I_h(x)$, 并估计误差.

19. 求 $f(x) = x^4$ 在 $[a, b]$ 上的分段埃尔米特插值, 并估计误差.

20. 给定数据表如下:

x_j	0.25	0.30	0.39	0.45	0.53
y_j	0.5000	0.5477	0.6245	0.6708	0.7280

试求三次样条插值 $S(x)$, 并满足条件:

$$\text{i) } S(0.25) = 1.0000, \quad S(0.53) = 0.6868;$$

$$\text{ii) } S(0.25) = S(0.53) = 0.$$

21. 若 $f(x) \in C^2[a, b]$, $S(x)$ 是三次样条函数, 证明:

$$\begin{aligned} \text{i) } & \int_a^b [f(x)]^2 dx - \int_a^b [S(x)]^2 dx \\ &= \int_a^b [f(x) - S(x)]^2 dx \\ &+ 2 \int_a^b S(x) [f(x) - S(x)] dx; \end{aligned}$$

ii) 若 $f(x_i) = S(x_i)$ ($i=0, 1, \dots, n$), 式中 x_i 为插值节点, 且 $a = x_0 < x_1 < \dots < x_n = b$, 则

$$\begin{aligned} & \int_a^b S(x) [f(x) - S(x)] dx \\ &= S(b)[f(b) - S(b)] - S(a)[f(a) - S(a)]. \end{aligned}$$

第3章 函数逼近与曲线拟合

3.1 函数逼近的基本概念

3.1.1 函数逼近与函数空间

在数值计算中经常要计算函数值,如计算机中计算基本初等函数及其他特殊函数;当函数只在有限点集上给定函数值,要在包含该点集的区间上用公式给出函数的简单表达式,这些都涉及到在区间 $[a, b]$ 上用简单函数逼近已知复杂函数的问题,这就是函数逼近问题.上章讨论的插值法就是函数逼近问题的一种.本章讨论的函数逼近,是指“对函数类 A 中给定的函数 $f(x)$,记作 $f(x)$

A ,要求在另一类简单的便于计算的函数类 B 中求函数 $p(x)$ B ,使 $p(x)$ 与 $f(x)$ 的误差在某种度量意义下最小”.函数类 A 通常是区间 $[a, b]$ 上的连续函数,记作 $C[a, b]$,称为连续函数空间,而函数类 B 通常为 n 次多项式,有理函数或分段低次多项式等.函数逼近是数值分析的基础,为了在数学上描述更精确,先要介绍代数和分析中一些基本概念及预备知识.

数学上常把在各种集合中引入某些不同的确定关系称为赋予集合以某种空间结构,并将这样的集合称为空间.例如将所有实 n 维向量组成的集合,按向量加法及向量与数的乘法构成实数域上的线性空间,记作 \mathbf{R}^n ,称为 n 维向量空间.类似地,对次数不超过 n (n 为正整数)的实系数多项式全体,按通常多项式与多项式加法及数与多项式乘法也构成数域 \mathbf{R} 上一个线性空间,用 H_n 表示,称为多项式空间.所有定义在 $[a, b]$ 上的连续函数集合,按函数加法

和数与函数乘法构成数域 \mathbf{R} 上的线性空间, 记作 $C[a, b]$. 类似地记 $C^p[a, b]$ 为具有 p 阶连续导数的函数空间.

定义 1 设集合 S 是数域 \mathbf{P} 上的线性空间, 元素 x_1, \dots, x_n 在 S , 如果存在不全为零的数 $a_1, \dots, a_n \in \mathbf{P}$, 使得

$$a_1 x_1 + \dots + a_n x_n = 0, \quad (1.1)$$

则称 x_1, \dots, x_n 线性相关. 否则, 若等式(1.1)只对 $a_1 = a_2 = \dots = a_n = 0$ 成立, 则称 x_1, \dots, x_n 线性无关.

若线性空间 S 是由 n 个线性无关元素 x_1, \dots, x_n 生成的, 即对 $\forall x \in S$ 都有

$$x = a_1 x_1 + \dots + a_n x_n,$$

则 x_1, \dots, x_n 称为空间 S 的一组基, 记为 $S = \text{span}\{x_1, \dots, x_n\}$, 并称空间 S 为 n 维空间, 系数 a_1, \dots, a_n 称为 x 在基 x_1, \dots, x_n 下的坐标, 记作 (a_1, \dots, a_n) , 如果 S 中有无限个线性无关元素 x_1, \dots, x_n, \dots , 则称 S 为无限维线性空间.

下面考察次数不超过 n 次的多项式集合 H_n , 其元素 $p(x)$ 在 H_n 表示为

$$p(x) = a_0 + a_1 x + \dots + a_n x^n, \quad (1.2)$$

它由 $n+1$ 个系数 (a_0, a_1, \dots, a_n) 唯一确定. $1, x, \dots, x^n$ 线性无关, 它是 H_n 的一组基, 故 $H_n = \text{span}\{1, x, \dots, x^n\}$, 且 (a_0, a_1, \dots, a_n) 是 $p(x)$ 的坐标向量, H_n 是 $n+1$ 维的.

对连续函数 $f(x) \in C[a, b]$, 它不能用有限个线性无关的函数表示, 故 $C[a, b]$ 是无限维的, 但它的任一元素 $f(x) \in C[a, b]$ 均可用有限维的 $p(x) \in H_n$ 逼近, 使误差 $\max_{a \leq x \leq b} |f(x) - p(x)| < \epsilon$ (ϵ 为任给的小正数), 这就是著名的魏尔斯特拉斯(Weierstrass)定理.

定理 1 设 $f(x) \in C[a, b]$, 则对任何 $\epsilon > 0$, 总存在一个代数多项式 $p(x)$, 使

$$|f(x) - p(x)| <$$

在 $[a, b]$ 上一致成立.

这定理已在“数学分析”中证明过。这里需要说明的是在许多证明方法中，伯恩斯坦（）1912年给出的证明是一种构造性证明。他根据函数整体逼近的特性构造出伯恩斯坦多项式

$$B_n(f, x) = \sum_{k=0}^n f\left(\frac{k}{n}\right) P_k(x); \quad (1.3)$$

其中

$$P_k(x) = \frac{n}{k} x^k (1-x)^{n-k},$$

$\frac{n}{k} = \frac{n(n-1)\dots(n-k+1)}{k!}$ 为二项式展开系数，并证明了 $\lim_n B_n(f, x) = f(x)$ 在 $[0, 1]$ 上一致成立；若 $f(x)$ 在 $[0, 1]$ 上 m 阶导数连续，则

$$\lim_n B_n^{(m)}(f, x) = f^{(m)}(x).$$

这不但证明了定理 1，而且由 (1.3) 给出了 $f(x)$ 的一个逼近多项式。它与拉格朗日插值多项式

$$L_n(x) = \sum_{k=0}^n f(x_k) l_k(x), \quad l_k(x) = 1$$

很相似，对 $B_n(f, x)$ ，当 $f(x) = 1$ 时也有关系式

$$\sum_{k=0}^n P_k(x) = \sum_{k=0}^n \frac{n}{k} x^k (1-x)^{n-k} = 1. \quad (1.4)$$

这只要在恒等式

$$(x+y)^n = \sum_{k=0}^n \frac{n}{k} x^k y^{n-k}$$

中令 $y = 1 - x$ 就可得到。但这里当 $x \in [0, 1]$ 时还有 $P_k(x) = 0$ ，于是

$$\sum_{k=0}^n |P_k(x)| = \sum_{k=0}^n P_k(x) = 1$$

是有界的，因而只要 $|f(x)|$ 对任意 $x \in [0, 1]$ 成立，则

$$|B_n(f, x)| \leq \max_{x \in [a, b]} |f(x)| \sum_{k=0}^n |P_k(x)|$$

有界, 故 $B_n(f, x)$ 是稳定的. 至于拉格朗日插值多项式 $L_n(x)$, 由

于 $\sum_{k=0}^n |l_k(x)|$ 无界, 因而不能保证高阶插值的稳定性与收敛性.

相比之下, 多项式 $B_n(f, x)$ 有良好的逼近性质, 但它收敛太慢, 比三次样条逼近效果差得多, 实际中很少被使用.

更一般地, 可用一组在 $C[a, b]$ 上线性无关的函数集合

$\{e_i(x)\}_{i=0}^n$ 来逼近 $f(x) \in C[a, b]$, 元素 $e(x) = \text{span}\{e_0(x),$

$e_1(x), \dots, e_n(x)\}$ $\subset C[a, b]$, 表示为

$$e(x) = a_0 e_0(x) + a_1 e_1(x) + \dots + a_n e_n(x). \quad (1.5)$$

函数逼近问题就是对任何 $f \in C[a, b]$, 在子空间 S 中找一个元素 $e^*(x)$, 使 $|f(x) - e^*(x)|$ 在某种意义上最小.

3.1.2 范数与赋范线性空间

为了对线性空间中元素大小进行衡量, 需要引进范数定义, 它是 \mathbf{R}^n 空间中向量长度概念的直接推广.

定义 2 设 S 为线性空间, $x \in S$, 若存在唯一实数 $\|\cdot\|$, 满足条件:

(1) $\|x\| \geq 0$, 当且仅当 $x = 0$ 时, $\|x\| = 0$; (正定性)

(2) $\|x\| = \|y\| \iff x = y$, $x, y \in \mathbf{R}^n$; (齐次性)

(3) $\|x + y\| \leq \|x\| + \|y\|$, $x, y \in S$. (三角不等式)

则称 $\|\cdot\|$ 为线性空间 S 上的范数, S 与 $\|\cdot\|$ 一起称为赋范线性空间, 记为 X .

例如, 在 \mathbf{R}^n 上的向量 $\mathbf{x} = (x_1, \dots, x_n)^T \in \mathbf{R}^n$, 三种常用范数为

$$\|\mathbf{x}\|_1 = \max_{1 \leq i \leq n} |x_i|, \quad \text{称为 } 1\text{-范数或最大范数},$$

$$x_1 = \left(\sum_{i=1}^n |x_i| \right)^{\frac{1}{n}}, \quad \text{称为 1-范数,}$$

$$x_2 = \left(\sum_{i=1}^n x_i^2 \right)^{\frac{1}{2}}, \quad \text{称为 2-范数.}$$

类似地对连续函数空间 $C[a, b]$, 若 $f \in C[a, b]$ 可定义三种常用范数如下:

$$\|f\|_\infty = \max_{a \leq x \leq b} |f(x)|, \quad \text{称为 } \infty\text{-范数,}$$

$$\|f\|_1 = \int_a^b |f(x)| dx, \quad \text{称为 1-范数,}$$

$$\|f\|_2 = \left(\int_a^b f^2(x) dx \right)^{\frac{1}{2}}, \quad \text{称为 2-范数.}$$

可以验证这样定义的范数均满足定义 2 中的三个条件.

3.1.3 内积与内积空间

在线性代数中, \mathbf{R}^n 中两个向量 $\mathbf{x} = (x_1, \dots, x_n)^T$ 及 $\mathbf{y} = (y_1, \dots, y_n)^T$ 的内积定义为

$$(x, y) = x_1 y_1 + \dots + x_n y_n.$$

若将它推广到一般的线性空间 X , 则有下面的定义.

定义 3 设 X 是数域 \mathbf{K} (\mathbf{R} 或 \mathbf{C}) 上的线性空间, 对 " $u, v \in X$, 有 \mathbf{K} 中一个数与之对应, 记为 (u, v) , 它满足以下条件:

$$(1) (u, v) = (\overline{v}, u), \quad " u, v \in X;$$

$$(2) (u, v) = (u, v), \quad \mathbf{K}, u, v \in X;$$

$$(3) (u+v, w) = (u, w) + (v, w), \quad " u, v, w \in X;$$

$$(4) (u, u) \geq 0, \text{ 当且仅当 } u=0 \text{ 时, } (u, u)=0.$$

则称 (u, v) 为 X 上 u 与 v 的内积. 定义了内积的线性空间称为内积空间. 定义中(1)的右端 (\overline{u}, v) 称为 (u, v) 的共轭, 当 \mathbf{K} 为实数域 \mathbf{R} 时 $(u, v) = (v, u)$.

如果 $(u, v) = 0$, 则称 u 与 v 正交, 这是向量相互垂直概念的推

广. 关于内积空间性质有以下重要定理.

定理2 设 X 为一个内积空间, 对 " $u, v \in X$, 有

$$|\langle u, v \rangle|^2 = (\langle u, u \rangle)(\langle v, v \rangle). \quad (1.6)$$

称为柯西-施瓦茨(Cauchy-Schwarz)不等式.

证明 当 $v=0$ 时(1.6)式显然成立. 现设 $v \neq 0$, 则 $\langle v, v \rangle > 0$, 且对任何数 有

$$0 = \langle u + v, u + v \rangle = \langle u, u \rangle + 2 \langle u, v \rangle + \langle v, v \rangle.$$

取 $\lambda = -\langle u, v \rangle / \langle v, v \rangle$, 代入上式右端, 得

$$\langle u, u \rangle - 2 \frac{\langle u, v \rangle}{\langle v, v \rangle} + \frac{\langle u, v \rangle}{\langle v, v \rangle}^2 \geq 0,$$

即得 $v \neq 0$ 时

$$|\langle u, v \rangle|^2 \leq (\langle u, u \rangle)(\langle v, v \rangle). \quad \text{证毕.}$$

定理3 设 X 为一个内积空间, $u_1, u_2, \dots, u_n \in X$, 矩阵

$$\mathbf{G} = \begin{pmatrix} \langle u_1, u_1 \rangle & \langle u_1, u_2 \rangle & \dots & \langle u_1, u_n \rangle \\ \langle u_2, u_1 \rangle & \langle u_2, u_2 \rangle & \dots & \langle u_2, u_n \rangle \\ \dots & \dots & \dots & \dots \\ \langle u_n, u_1 \rangle & \langle u_n, u_2 \rangle & \dots & \langle u_n, u_n \rangle \end{pmatrix} \quad (1.7)$$

称为格拉姆(Gram)矩阵, 则 \mathbf{G} 非奇异的充分必要条件是 u_1, u_2, \dots, u_n 线性无关.

证明 \mathbf{G} 非奇异等价于 $\det \mathbf{G} \neq 0$, 其充要条件是齐次方程组

$$\sum_{j=1}^n \langle u_j, u_k \rangle = \sum_{j=1}^n \langle u_j, u_k \rangle_j = 0, \quad k = 1, 2, \dots, n \quad (1.8)$$

只有零解; 而

$$\begin{aligned} \sum_{j=1}^n \langle u_j, u_1 \rangle + \sum_{j=1}^n \langle u_j, u_2 \rangle + \dots + \sum_{j=1}^n \langle u_j, u_n \rangle &= 0 \\ \sum_{j=1}^n \langle u_j, u_1 \rangle, \sum_{j=1}^n \langle u_j, u_2 \rangle &= 0 \\ \sum_{j=1}^n \langle u_j, u_k \rangle, \quad k = 1, 2, \dots, n. \end{aligned} \quad (1.9)$$

从以上等价关系可知, $\det \mathbf{G} = 0$ 等价于从(1.8)推出 $u_1 = u_2 = \dots = u_n = 0$, 而后者等价于从(1.9)推出 $u_1 = u_2 = \dots = u_n = 0$, 即 u_1, u_2, \dots, u_n 线性无关. 证毕.

在内积空间 X 上可以由内积导出一种范数, 即对于 $u \in X$, 记

$$\|u\| = (\langle u, u \rangle)^{\frac{1}{2}}, \quad (1.10)$$

容易验证它满足范数定义的三条性质, 其中三角不等式

$$\|u + v\| \leq \|u\| + \|v\| \quad (1.11)$$

可由定理2直接得出, 即

$$\begin{aligned} (\|u + v\|)^2 &= \|u\|^2 + 2\langle u, v \rangle + \|v\|^2 \\ &= (\langle u, u \rangle + 2\langle u, v \rangle + \langle v, v \rangle)^{\frac{1}{2}} \\ &= (\langle u + v, u + v \rangle)^{\frac{1}{2}} = \|u + v\|^2, \end{aligned}$$

两端开方即得(1.11).

例1 \mathbf{R}^n 与 \mathbf{C}^n 的内积. 设 $x, y \in \mathbf{R}^n$, $\mathbf{x} = (x_1, \dots, x_n)^T$, $\mathbf{y} = (y_1, \dots, y_n)^T$, 则其内积定义为

$$\langle x, y \rangle = \sum_{i=1}^n x_i y_i. \quad (1.12)$$

由此导出的向量2-范数为

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

若给定实数 $\alpha_i > 0$ ($i = 1, \dots, n$), 称 $\{\alpha_i\}$ 为权系数, 则在 \mathbf{R}^n 上可定义加权内积为

$$\langle x, y \rangle = \sum_{i=1}^n \alpha_i x_i y_i, \quad (1.13)$$

相应的范数为

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

不难验证(1.13)给出的 $\langle x, y \rangle$ 满足内积定义的4条性质. 当 $\alpha_i = 1$ ($i = 1, \dots, n$) 时, (1.13) 就是(1.12).

如果 $x, y \in \mathbf{C}^n$, 带权内积定义为

$$(x, y) = \sum_{i=1}^n x_i y_i \bar{\omega}_i, \quad (1.14)$$

这里 $\{\omega_i\}$ 仍为正实数序列, $\bar{\omega}_i$ 为 ω_i 的共轭.

在 $C[a, b]$ 上也可以类似定义带权内积, 为此先给出权函数的定义.

定义4 设 $[a, b]$ 是有限或无限区间, 在 $[a, b]$ 上的非负函数 (x) 满足条件:

$$(1) \int_a^b x^k (x) dx \text{ 存在且为有限值} (k=0, 1, \dots),$$

(2) 对 $[a, b]$ 上的非负连续函数 $g(x)$, 如果 $\int_a^b g(x) (x) dx = 0$, 则 $g(x) = 0$.

则称 (x) 为 $[a, b]$ 上的一个权函数.

例2 $C[a, b]$ 上的内积. 设 $f(x), g(x) \in C[a, b]$, (x) 是 $[a, b]$ 上给定的权函数, 则可定义内积

$$(f(x), g(x)) = \int_a^b f(x) g(x) (x) dx. \quad (1.15)$$

容易验证它满足内积定义的4条性质, 由此内积导出的范数为

$$\|f(x)\|_2 = (f(x), f(x))^{\frac{1}{2}} = \left(\int_a^b f(x)^2 (x) dx \right)^{\frac{1}{2}} \quad (1.16)$$

称(1.15)和(1.16)为带权 (x) 的内积和范数, 特别常用的是 $(x)_1$ 的情形, 即

$$(f(x), g(x)) = \int_a^b f(x) g(x) (x) dx,$$

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

若 $\phi_0, \phi_1, \dots, \phi_n$ 是 $C[a, b]$ 中的线性无关函数族, 记 $\Phi = \text{span}\{\phi_0, \phi_1, \dots, \phi_n\}$, 它的格拉姆矩阵为

$$\begin{aligned}
 & \quad (0, 0) \quad (0, 1) \quad \dots \quad (0, n) \\
 \mathbf{G} = \mathbf{G}(0, 1, \dots, n) = & \quad (1, 0) \quad (1, 1) \quad \dots \quad (1, n) \\
 & \quad \dots \quad \dots \quad \dots \\
 & \quad (n, 0) \quad (n, 1) \quad \dots \quad (n, n)
 \end{aligned} \tag{1.17}$$

根据定理 3 可知 $(0, 1, \dots, n)$ 线性无关的充要条件是 $\det \mathbf{G}(0, 1, \dots, n) \neq 0$.

3.2 正交多项式

正交多项式是函数逼近的重要工具，在数值积分中也有重要应用.

3.2.1 正交函数族与正交多项式

定义 5 若 $f(x), g(x) \in C[a, b]$, (x) 为 $[a, b]$ 上的权函数且满足

$$(f(x), g(x)) = \int_a^b (x) f(x) g(x) dx = 0, \tag{2.1}$$

则称 $f(x)$ 与 $g(x)$ 在 $[a, b]$ 上带权 (x) 正交. 若函数族 $_0(x), _1(x), \dots, _n(x), \dots$ 满足关系

$$(_j, _k) = \int_a^b (x) _j(x) _k(x) dx = \begin{cases} 0, & j \neq k, \\ A_k > 0, & j = k. \end{cases} \tag{2.2}$$

则称 $\{_k(x)\}$ 是 $[a, b]$ 上带权 (x) 的正交函数族；若 $A_k = 1$ ，则称之为标准正交函数族.

例如，三角函数族

$$1, \cos x, \sin x, \cos 2x, \sin 2x, \dots$$

就是在区间 $[-\pi, \pi]$ 上的正交函数族. 因为对 $k=1, 2, \dots$ 有

$$(1, 1) = 2, \quad (\sin kx, \sin kx) = (\cos kx, \cos kx) = ,$$

而对 $k, j = 1, 2, \dots$, 当 $k = j$ 时有 $(\cos kx, \sin kx) = (1, \cos kx) = (1, \sin kx) = 0$;

$$(\cos kx, \cos jx) = (\sin kx, \sin jx) = (\cos kx, \sin jx) = 0.$$

定义 6 设 $\{n(x)\}$ 是 $[a, b]$ 上首项系数 $a_n \neq 0$ 的 n 次多项式, (x) 为 $[a, b]$ 上权函数, 如果多项式序列 $\{n(x)\}_0^\infty$ 满足关系式 (2.2), 则称多项式序列 $\{n(x)\}_0^\infty$ 为在 $[a, b]$ 上带权 (x) 正交, 称 $\{n(x)\}$ 为 $[a, b]$ 上带权 (x) 的 n 次正交多项式.

只要给定区间 $[a, b]$ 及权函数 (x) , 均可由一族线性无关的幂函数 $\{1, x, \dots, x^n, \dots\}$, 利用逐个正交化手续构造出正交多项式序列 $\{n(x)\}_0^\infty$:

$$\begin{aligned} n_0(x) &= 1, \\ n_n(x) &= x^n - \sum_{j=0}^{n-1} \frac{(x^n, j(x))}{(j(x), j(x))} j(x) \\ (n &= 1, 2, \dots). \end{aligned} \quad (2.3)$$

这样得到的正交多项式序列有以下性质:

(1) $n(x)$ 是具有最高次项系数为 1 的 n 次多项式.

(2) 任何 n 次多项式 $P_n(x)$ H_n 均可表示为 $n_0(x), n_1(x), \dots, n_n(x)$ 的线性组合.

(3) 当 $k = j$ 时, $(j(x), k(x)) = 0$, 且 $n_k(x)$ 与任一次数小于 k 的多项式正交.

(4) 成立递推关系

$$n_{n+1}(x) = (x - n_n(x)) n_n(x) - n_{n-n-1}(x) \quad (n = 0, 1, \dots), \quad (2.4)$$

其中

$$n_0(x) = 1, \quad n_{-1}(x) = 0,$$

$$n_n = (x - n_n(x), n_n(x)) / (n_n(x), n_n(x)),$$

$$n_n = (n_n(x), n_n(x)) / (n_{n-1}(x), n_{n-1}(x)) \quad (n = 1, 2, \dots),$$

这里 $(x_n(x), x_m(x)) = \int_a^b x_n^2(x) x_m(x) dx$.

(5) 设 $\{x_n(x)\}_0$ 是在 $[a, b]$ 上带权 $w(x)$ 的正交多项式序列, 则 $x_n(x) (n \geq 1)$ 的 n 个根都是在区间 (a, b) 内的单重实根.

以上性质证明可见 [2], 下面给出几种常见的正交多项式.

3.2.2 勒让德多项式

当区间为 $[-1, 1]$, 权函数 $w(x) = 1$ 时, 由 $\{1, x, \dots, x^n, \dots\}$ 正交化得到的多项式就称为勒让德(Legendre)多项式, 并用 $P_0(x)$, $P_1(x), \dots, P_n(x), \dots$ 表示. 这是勒让德于 1785 年引进的. 1814 年罗德利克(Rodrigul)给出了简单的表达式

$$P_0(x) = 1, \quad P_n(x) = \frac{1}{2^n n!} \frac{d^n}{dx^n} \{(x^2 - 1)^n\} \quad (n = 1, 2, \dots), \quad (2.5)$$

由于 $(x^2 - 1)^n$ 是 $2n$ 次多项式, 求 n 阶导数后得

$$P_n(x) = \frac{1}{2^n n!} (2n)(2n-1)\dots(n+1)x^n + a_{n-1}x^{n-1} + \dots + a,$$

于是得首项 x^n 的系数 $a_n = \frac{(2n)!}{2^n (n!)^2}$. 显然最高项系数为 1 的勒让德多项式为

$$P_n(x) = \frac{n!}{(2n)!} \frac{d^n}{dx^n} \{(x^2 - 1)^n\}. \quad (2.6)$$

勒让德多项式有下述几个重要性质:

性质 1 正交性

$$\int_{-1}^1 P_n(x) P_m(x) dx = \begin{cases} 0, & m \neq n; \\ \frac{2}{2n+1}, & m = n. \end{cases} \quad (2.7)$$

证明 令 $w(x) = (x^2 - 1)^n$, 则

$$w^{(k)}(\pm 1) = 0 \quad (k = 0, 1, \dots, n-1).$$

设 $Q(x)$ 是在区间 $[-1, 1]$ 上有 n 阶连续可微的函数, 由分部积分知

$$\begin{aligned} \int_{-1}^1 P_n(x) Q(x) dx &= \frac{1}{2^n n!} \int_{-1}^1 Q(x)^{(n)}(x) dx \\ &= -\frac{1}{2^n n!} \int_{-1}^1 Q^{(n)}(x)(x) dx \\ &= \dots \\ &= \frac{(-1)^n}{2^n n!} \int_{-1}^1 Q^{(n)}(x) dx. \end{aligned}$$

下面分两种情况讨论.

(1) 若 $Q(x)$ 是次数小于 n 的多项式, 则 $Q^{(n)}(x) = 0$, 故得

$$\int_{-1}^1 P_m(x) P_n(x) dx = 0, \quad \text{当 } n \neq m.$$

(2) 若 $Q(x) = P_n(x) = \frac{1}{2^n n!} (x)^{(n)}(x) = \frac{(2n)!}{2^n (n!)^2} x^n + \dots$,

$$Q^{(n)}(x) = P_n^{(n)}(x) = \frac{(2n)!}{2^n n!},$$

于是

$$\begin{aligned} \int_{-1}^1 P_n^2(x) dx &= \frac{(-1)^n (2n)!}{2^{2n} (n!)^2} \int_{-1}^1 (x^2 - 1)^n dx \\ &= \frac{(2n)!}{2^{2n} (n!)^2} \int_{-1}^1 (1 - x^2)^n dx. \end{aligned}$$

由于

$$\int_0^1 (1 - x^2)^n dx = \int_0^{\frac{\pi}{2}} \cos^{2n+1} t dt = \frac{2 \cdot 4 \cdot \dots \cdot (2n)}{1 \cdot 3 \cdot \dots \cdot (2n+1)},$$

故

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

于是(2.7)得证.

性质 2 奇偶性

$$P_n(-x) = (-1)^n P_n(x). \quad (2.8)$$

由于 $(x) = (x^2 - 1)^n$ 是偶次多项式, 经过偶次求导仍为偶次多项式, 经过奇次求导则为奇次多项式, 故 n 为偶数时 $P_n(x)$ 为偶函数, n 为奇数时 $P_n(x)$ 为奇函数, 于是(2.8)成立.

性质3 递推关系

考虑 $n+1$ 次多项式 $xP_n(x)$, 它可表示为

$$xP_n(x) = a_0 P_0(x) + a_1 P_1(x) + \dots + a_{n+1} P_{n+1}(x).$$

两边乘 $P_k(x)$, 并从 -1 到 1 积分, 得

$$\int_{-1}^1 xP_n(x)P_k(x)dx = a_k \int_{-1}^1 P_k^2(x)dx.$$

当 $k = n-2$ 时, $xP_k(x)$ 次数小于等于 $n-1$, 上式左端积分为 0, 故得 $a_k = 0$. 当 $k = n$ 时, $xP_n^2(x)$ 为奇函数, 左端积分仍为 0, 故 $a_n = 0$. 于是

$$xP_n(x) = a_{n-1} P_{n-1}(x) + a_{n+1} P_{n+1}(x),$$

其中

$$\begin{aligned} a_{n-1} &= \frac{2n-1}{2} \int_{-1}^1 xP_n(x)P_{n-1}(x)dx \\ &= \frac{2n-1}{2} \cdot \frac{2n}{4n^2-1} = \frac{n}{2n+1}, \\ a_{n+1} &= \frac{2n+3}{2} \int_{-1}^1 xP_n(x)P_{n+1}(x)dx \\ &= \frac{2n+3}{2} \cdot \frac{2(n+1)}{(2n+1)(2n+3)} = \frac{n+1}{2n+1}, \end{aligned}$$

从而得到以下的递推公式

$$(n+1)P_{n+1}(x) = (2n+1)xP_n(x) - nP_{n-1}(x) \quad (n=1, 2, \dots), \quad (2.9)$$

由 $P_0(x) = 1$, $P_1(x) = x$, 利用(2.9)就可推出

$$P_2(x) = (3x^2 - 1)/2,$$

$$P_3(x) = (5x^3 - 3x)/2,$$

$$P_4(x) = (35x^4 - 30x^2 + 3)/8,$$

$$P_5(x) = (63x^5 - 70x^3 + 15x)/8,$$

$$P_6(x) = (231x^6 - 315x^4 + 105x^2 - 5)/16,$$

...

图 3-1 给出了 $P_0(x), P_1(x), P_2(x), P_3(x)$ 的图形 .

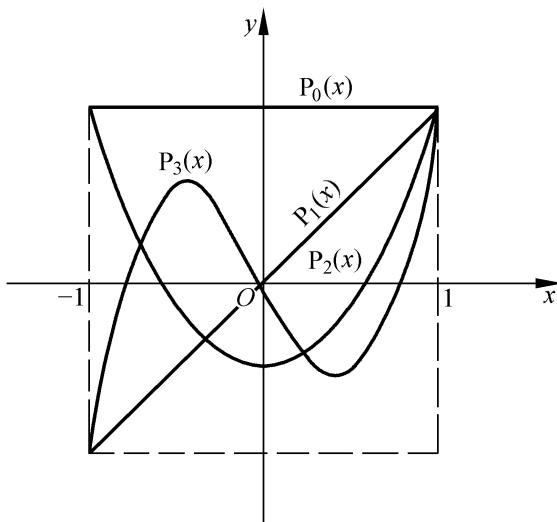


图 3-1

性质 4 $P_n(x)$ 在区间 $[-1, 1]$ 内有 n 个不同的实零点 .

3.2.3 切比雪夫多项式

当权函数 $\omega(x) = \frac{1}{1-x^2}$, 区间为 $[-1, 1]$ 时, 由序列 $\{1, x, \dots, x^n, \dots\}$ 正交化得到的正交多项式就是切比雪夫 (Chebyshev) 多项式, 它可表示为

$T_n(x) = \cos(n \arccos x), \quad |x| \leq 1. \quad (2.10)$

若令 $x = \cos \theta$, 则 $T_n(x) = \cos n\theta$.

切比雪夫多项式有很多重要性质:

性质 5 递推关系

$$T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x) \quad (n = 1, 2, \dots),$$

$$T_0(x) = 1, \quad T_1(x) = x. \quad (2.11)$$

这只要由三角恒等式

$$\cos(n+1) = 2\cos \cos n - \cos(n-1) \quad (n \geq 1)$$

令 $x = \cos$ 即得. 由(2.11)就可推出

$$T_2(x) = 2x^2 - 1,$$

$$T_3(x) = 4x^3 - 3x,$$

$$T_4(x) = 8x^4 - 8x^2 + 1,$$

$$T_5(x) = 16x^5 - 20x^3 + 5x,$$

$$T_6(x) = 32x^6 - 48x^4 + 18x^2 - 1,$$

...

$T_0(x), T_1(x), T_2(x), T_3(x)$ 的函数图形见图 3-2.

由递推关系(2.11)还可得到 $T_n(x)$ 的最高次项系数是 2^{n-1} ($n \geq 1$).

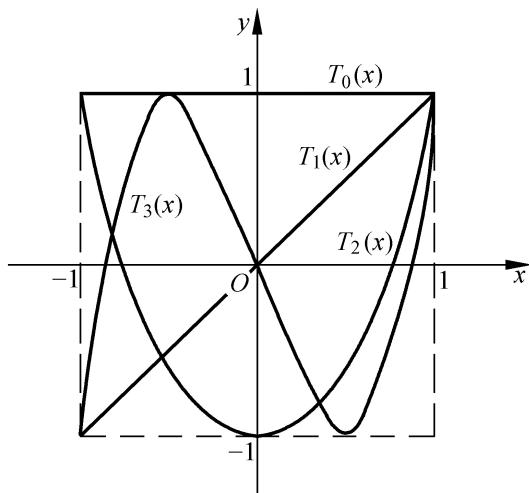


图 3-2

性质 6 切比雪夫多项式 $\{T_k(x)\}$ 在区间 $[-1, 1]$ 上带权 $w(x) = 1/\sqrt{1-x^2}$ 正交, 且

$$0, \quad n \neq m;$$

$$\int_{-1}^1 \frac{T_n(x) T_m(x) dx}{\sqrt{1-x^2}} = \begin{cases} 0, & n \neq m \\ \frac{\pi}{2}, & n = m = 0 \\ , & n = m \neq 0 \end{cases} \quad (2.12)$$

事实上,令 $x = \cos d$, 则 $dx = -\sin d$, 于是

$$0, \quad n = m;$$

$$\begin{aligned} \frac{T_n(x) T_m(x)}{1 - x^2} dx &= \int_0^{\pi} \cos n \cos m d = \frac{1}{2}, \quad n = m = 0; \\ &\quad , \quad n = m = 0. \end{aligned}$$

性质7 $T_{2k}(x)$ 只含 x 的偶次幂, $T_{2k+1}(x)$ 只含 x 的奇次幂.

这性质由递推关系直接得到.

性质8 $T_n(x)$ 在区间 $[-1, 1]$ 上有 n 个零点

$$x_k = \cos \frac{2k-1}{2n}, \quad k = 1, \dots, n.$$

此外, 实际计算中时常要求 x^n 用 $T_0(x), T_1(x), \dots, T_n(x)$ 的线性组合表示, 其公式为

$$x^n = \sum_{k=0}^{\frac{n}{2}} \binom{n}{k} T_{n-2k}(x). \quad (2.13)$$

这里规定 $T_0(x) = 1$. $n = 1, 2, \dots, 6$ 时的结果如下:

$$1 = T_0(x),$$

$$x = T_1(x),$$

$$x^2 = \frac{1}{2}(T_0(x) + T_2(x)),$$

$$x^3 = \frac{1}{4}(3T_1(x) + T_3(x)),$$

$$x^4 = \frac{1}{8}(3T_0(x) + 4T_2(x) + T_4(x)),$$

$$x^5 = \frac{1}{16}(10T_1(x) + 5T_3(x) + T_5(x)),$$

$$x^6 = \frac{1}{32}(10T_0(x) + 15T_2(x) + 6T_4(x) + T_6(x)).$$

3.2.4 其他常用的正交多项式

一般说, 如果区间 $[a, b]$ 及权函数 (x) 不同, 则得到的正交多项式也不同. 除上述两种最重要的正交多项式外, 下面再给出三种较常用的正交多项式.

1. 第二类切比雪夫多项式

在区间 $[-1, 1]$ 上带权 $(x) = 1 - x^2$ 的正交多项式称为第二类切比雪夫多项式, 其表达式为

$$U_n(x) = \frac{\sin[(n+1)\arccos x]}{1-x^2}. \quad (2.14)$$

令 $x = \cos \theta$, 可得

$$\begin{aligned} \int_{-1}^1 U_n(x) U_m(x) (1-x^2) dx &= \int_0^\pi \sin(n+1)\theta \sin(m+1)\theta d\theta \\ &= \frac{0, m \neq n;}{2, m = n}, \end{aligned}$$

即 $\{U_n(x)\}$ 是 $[-1, 1]$ 上带权 $1 - x^2$ 的正交多项式族. 还可得到递推关系式

$$\begin{aligned} U_0(x) &= 1, & U_1(x) &= 2x, \\ U_{n+1}(x) &= 2xU_n(x) - U_{n-1}(x) \quad (n = 1, 2, \dots). \end{aligned}$$

2. 拉盖尔多项式

在区间 $[0, +\infty)$ 上带权 e^{-x} 的正交多项式称为拉盖尔(Laguerre)多项式, 其表达式为

$$L_n(x) = e^x \frac{d^n}{dx^n} (x^n e^{-x}). \quad (2.15)$$

它也具有正交性质

$$\int_0^\infty e^{-x} L_n(x) L_m(x) dx = \begin{cases} 0, & m \neq n; \\ (n!)^2, & m = n, \end{cases}$$

和递推关系

$$\begin{aligned} L_0(x) &= 1, & L_1(x) &= 1 - x, \\ L_{n+1}(x) &= (1 + 2n - x)L_n(x) - n^2 L_{n-1}(x) \quad (n = 1, 2, \dots). \end{aligned}$$

3. 埃尔米特多项式

在区间 $(-, +)$ 上带权 e^{-x^2} 的正交多项式称为埃尔米特多项式, 其表达式为

$$H_n(x) = (-1)^n e^{-x^2} \frac{d^n}{dx^n}(e^{-x^2}), \quad (2.16)$$

它满足正交关系

$$\int_{-\infty}^{+\infty} e^{-x^2} H_m(x) H_n(x) dx = \begin{cases} 0, & m \neq n; \\ 2^n n!, & m = n, \end{cases}$$

并有递推关系

$$\begin{aligned} H_0(x) &= 1, & H_1(x) &= 2x, \\ H_{n+1}(x) &= 2xH_n(x) - 2nH_{n-1}(x) \quad (n = 1, 2, \dots). \end{aligned}$$

3.3 最佳一致逼近多项式

3.3.1 基本概念及其理论

本节讨论 $f \in C[a, b]$, 在 $H_n = \text{span}\{1, x, \dots, x^n\}$ 中求多项式 $P_n^*(x)$, 使其误差

$$f - P_n^* = \max_{a \leq x \leq b} |f(x) - P_n^*(x)| = \min_{P_n \in H_n} |f - P_n|.$$

这就是通常所谓最佳一致逼近或切比雪夫逼近问题. 为了说明这一概念, 先给出以下定义.

定义 7 设 $P_n(x) \in H_n$, $f(x) \in C[a, b]$, 称

$$(f, P_n) = f - P_n = \max_{a \leq x \leq b} |f(x) - P_n(x)| \quad (3.1)$$

为 $f(x)$ 与 $P_n(x)$ 在 $[a, b]$ 上的偏差.

显然 $(f, P_n) \geq 0$, (f, P_n) 的全体组成一个集合, 记为

$\{ (f, P_n) \}$, 它有下界 0. 若记集合的下确界为

$$E_n = \inf_{P_n \in H_n} \{ (f, P_n) \} = \inf_{P_n \in H_n} \max_{\substack{a \\ x \\ b}} |f(x) - P_n(x)|, \quad (3.2)$$

则称之为 $f(x)$ 在 $[a, b]$ 上的最小偏差.

定义 8 假定 $f(x) \in C[a, b]$, 若存在 $P_n^*(x) \in H_n$ 使得

$$(f, P_n^*) = E_n, \quad (3.3)$$

则称 $P_n^*(x)$ 是 $f(x)$ 在 $[a, b]$ 上的最佳一致逼近多项式或最小偏差逼近多项式, 简称最佳逼近多项式.

注意, 定义并未说明最佳逼近多项式是否存在, 但可证明下面的存在定理.

定理 4 若 $f(x) \in C[a, b]$, 则总存在 $P_n^*(x) \in H_n$, 使

$$|f(x) - P_n^*(x)| = E_n.$$

证明略, 可参考 [2].

为了研究最佳逼近多项式的特性, 先引进偏差点的定义.

定义 9 设 $f(x) \in C[a, b]$, $P(x) \in H_n$, 若在 $x = x_0$ 上有

$$|P(x_0) - f(x_0)| = \max_{\substack{a \\ x \\ b}} |P(x) - f(x)| = \mu,$$

就称 x_0 是 $P(x)$ 的偏差点.

若 $P(x_0) - f(x_0) = \mu$, 称 x_0 为“正”偏差点.

若 $P(x_0) - f(x_0) = -\mu$, 称 x_0 为“负”偏差点.

由于函数 $P(x) - f(x)$ 在 $[a, b]$ 上连续, 因此, 至少存在一个点 $x_0 \in [a, b]$, 使

$$|P(x_0) - f(x_0)| = \mu,$$

也就是说 $P(x)$ 的偏差点总是存在的. 下面给出反映最佳逼近多项式特征的切比雪夫定理.

定理 5 $P(x) \in H_n$ 是 $f(x) \in C[a, b]$ 的最佳逼近多项式的充分必要条件是 $P(x)$ 在 $[a, b]$ 上至少有 $n+2$ 个轮流为“正”、“负”的偏差点, 即有 $n+2$ 个点 $a = x_1 < x_2 < \dots < x_{n+2} = b$, 使

$$P(x_k) - f(x_k) = (-1)^k \quad P(x) - f(x) \quad , \quad = \pm 1, \quad (3.4)$$

这样的点组称为切比雪夫交错点组.

证明 只证充分性. 假定在 $[a, b]$ 上有 $n+2$ 个点使(3.4)成立, 要证明 $P(x)$ 是 $f(x)$ 在 $[a, b]$ 上的最佳逼近多项式. 用反证法, 若存在 $Q(x) \in H_n$, $Q(x) \neq P(x)$, 使

$$f(x) - Q(x) < f(x) - P(x).$$

由于

$$P(x) - Q(x) = [P(x) - f(x)] - [Q(x) - f(x)]$$

在点 x_1, x_2, \dots, x_{n+2} 上的符号与 $P(x_k) - f(x_k)$ ($k=1, \dots, n+2$) 一致, 故 $P(x) - Q(x)$ 也在 $n+2$ 个点上轮流取“+”、“-”号. 由连续函数性质, 它在 $[a, b]$ 内有 $n+1$ 个零点, 但因 $P(x) - Q(x) \neq 0$ 是不超过 n 次的多项式, 它的零点不超过 n . 这矛盾说明假设不对, 故 $P(x)$ 就是所求最佳逼近多项式. 充分性得证. 必要性证明略, 可参看[5].

定理 5 说明用 $P(x)$ 逼近 $f(x)$ 的误差曲线 $y = P(x) - f(x)$ 是均匀分布的. 由这定理还可得以下重要推论.

推论 1 若 $f(x) \in C[a, b]$, 则在 H_n 中存在唯一的最佳逼近多项式.

证明略.

利用定理 5 可直接得到切比雪夫多项式 $T_n(x)$ 的一个重要性质, 即

定理 6 在区间 $[-1, 1]$ 上所有最高次项系数为 1 的 n 次多项式中, $\|T_n(x)\| = \frac{1}{2^{n-1}} \|T_n(x)\|$ 与零的偏差最小, 其偏差为 $\frac{1}{2^{n-1}}$.

证明 由于

$$T_n(x) = \frac{1}{2^{n-1}} T_n(x) = x^n - P_{n-1}^*(x),$$

$$\max_{[-1, 1]} |P_n(x)| = \frac{1}{2^{n-1}} \cdot \max_{[-1, 1]} |T_n(x)| = \frac{1}{2^{n-1}},$$

且点 $x_k = \cos \frac{k\pi}{n}$ ($k=0, 1, \dots, n$) 是 $T_n(x)$ 的切比雪夫交错点组, 由定理 5 可知, 区间 $[-1, 1]$ 上 x^n 在 H_{n-1} 中最佳逼近多项式为 $P_{n-1}^*(x)$, 即 $P_n(x)$ 是与零的偏差最小的多项式. 定理得证.

例 3 求 $f(x) = 2x^3 + x^2 + 2x - 1$ 在 $[-1, 1]$ 上的最佳 2 次逼近多项式.

解 由题意, 所求最佳逼近多项式 $P_2^*(x)$ 应满足

$$\max_{[-1, 1]} |f(x) - P_2^*(x)| = \min.$$

由定理 6 可知, 当

$$f(x) - P_2^*(x) = \frac{1}{2} T_3(x) = 2x^3 - \frac{3}{2}x$$

时, 多项式 $f(x) - P_2^*(x)$ 与零偏差最小, 故

$$\begin{aligned} P_2^*(x) &= f(x) - \frac{1}{2} T_3(x) \\ &= x^2 + \frac{7}{2}x - 1 \end{aligned}$$

就是 $f(x)$ 在 $[-1, 1]$ 上的最佳 2 次逼近多项式.

3.3.2 最佳一次逼近多项式

定理 5 给出了最佳逼近多项式 $P(x)$ 的特性, 但要求出 $P(x)$ 却相当困难. 下面讨论 $n=1$ 的情形. 假定 $f(x) \in C[a, b]$, 且 $f(x)$ 在 (a, b) 内不变号, 我们要求最佳一次逼近多项式 $P_1(x) = a + a_1 x$. 根据定理 5 可知至少有 3 个点 $a < x_1 < x_2 < x_3 < b$, 使

$$\begin{aligned} P_1(x_k) - f(x_k) &= (-1)^k \max_{a \leq x \leq b} |P_1(x) - f(x)| \\ &\quad (k = \pm 1, k = 1, 2, 3). \end{aligned}$$

由于 $f(x)$ 在 $[a, b]$ 上不变号, 故 $f(x)$ 单调, $f(x) - a_1$ 在

(a, b) 内只有一个零点, 记为 x_2 , 于是

$$P_1(x_2) - f(x_2) = a - f(x_2) = 0, \text{ 即 } f(x_2) = a.$$

另外两个偏差点必是区间端点, 即 $x_1 = a, x_2 = b$, 且满足

$$P_1(a) - f(a) = P_1(b) - f(b) = -[P_1(x_2) - f(x_2)].$$

由此得到

$$\begin{aligned} a_0 + a_1 a - f(a) &= a_0 + a_1 b - f(b); \\ a_0 + a_1 a - f(a) &= f(x_2) - (a_0 + a_1 x_2). \end{aligned} \quad (3.5)$$

解出

$$a_1 = \frac{f(b) - f(a)}{b - a} = f(x_2), \quad (3.6)$$

代入(3.5)得

$$a_0 = \frac{f(a) + f(x_2)}{2} - \frac{f(b) - f(a)}{b - a} \frac{a + x_2}{2}. \quad (3.7)$$

这就得到最佳一次逼近多项式 $P_1(x)$, 其几何意义如图 3-3 所示.

直线 $y = P_1(x)$ 与弦 MN 平行, 且通过 MQ 的中点 D , 其方程为

$$y = \frac{1}{2}[f(a) + f(x_2)] + a_1 x - \frac{a + x_2}{2}.$$

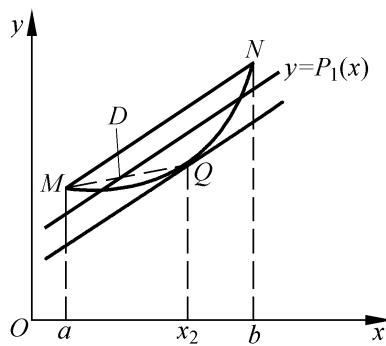


图 3-3

例 4 求 $f(x) = 1 + x^2$ 在 $[0, 1]$ 上的最佳一次逼近多项式.

解 由(3.6)可算出

$$a_1 = 2 - 1 = 0.414,$$

又 $f(x) = \frac{x}{1+x^2}$, 故 $\frac{x_2}{1+x_2^2} = 2 - 1$, 解得

$$x_2 = \frac{2 - 1}{2} = 0.4551, f(x_2) = 1 + x_2^2 = 1.0986.$$

由(3.7), 得

$$a_0 = \frac{1 + 1 + x_2^2}{2} = a_1 \frac{x_2}{2} = 0.955,$$

于是得 $1 + x^2$ 的最佳一次逼近多项式为

$$P_1(x) = 0.955 + 0.414x,$$

即 $1 + x^2 = 0.955 + 0.414x, 0 \leq x \leq 1$; (3.8)

误差限为

$$\max_{0 \leq x \leq 1} |1 + x^2 - P_1(x)| = 0.045.$$

在(3.8)中若令 $x = \frac{b}{a} - 1$, 则可得一个求根式的公式

$$a^2 + b^2 = 0.955a + 0.414b.$$

3.4 最佳平方逼近

3.4.1 最佳平方逼近及其计算

现在我们研究在区间 $[a, b]$ 上一般的最佳平方逼近问题. 对 $f(x) \in C[a, b]$ 及 $C[a, b]$ 中的一个子集 $S = \text{span}\{s_0(x), s_1(x), \dots, s_n(x)\}$, 若存在 $S^*(x)$, 使

$$\begin{aligned} \|f(x) - S^*(x)\|_2^2 &= \min_{S(x)} \|f(x) - S(x)\|_2^2 \\ &= \min_{S(x)} \int_a^b [f(x) - S(x)]^2 dx. \end{aligned} \quad (4.1)$$

则称 $S^*(x)$ 是 $f(x)$ 在子集 $C[a, b]$ 中的最佳平方逼近函数. 为了求 $S^*(x)$, 由(4.1)可知该问题等价于求多元函数

$$I(a_0, a_1, \dots, a_n) = \int_a^b [f(x) - \sum_{j=0}^n a_j x^j]^2 dx \quad (4.2)$$

的最小值. 由于 $I(a_0, a_1, \dots, a_n)$ 是关于 a_0, a_1, \dots, a_n 的二次函数, 利用多元函数求极值的必要条件

$$\frac{\partial I}{\partial a_k} = 0 \quad (k = 0, 1, \dots, n),$$

即

$$\frac{\partial I}{\partial a_k} = 2 \int_a^b [f(x) - \sum_{j=0}^n a_j x^j]^2 dx = 0 \quad (k = 0, 1, \dots, n),$$

于是有

$$\sum_{j=0}^n (\delta_{jk}(x), \delta_{jk}(x)) a_j = (f(x), \delta_{kk}(x)) \quad (k = 0, 1, \dots, n). \quad (4.3)$$

这是关于 a_0, a_1, \dots, a_n 的线性方程组, 称为法方程, 由于 $\delta_{00}(x), \delta_{11}(x), \dots, \delta_{nn}(x)$ 线性无关, 故系数 $\det \mathbf{G}(\delta_{00}, \delta_{11}, \dots, \delta_{nn}) \neq 0$, 于是方程组(4.3)有唯一解 $a_k^* (k = 0, 1, \dots, n)$, 从而得到

$$S^*(x) = a_0^* \delta_{00}(x) + \dots + a_n^* \delta_{nn}(x).$$

下面证明 $S^*(x)$ 满足(4.1), 即对任何 $S(x)$, 有

$$\int_a^b (f(x) - S^*(x))^2 dx = \int_a^b (f(x) - S(x))^2 dx. \quad (4.4)$$

为此只要考虑

$$\begin{aligned} D &= \int_a^b (f(x) - S(x))^2 dx - \int_a^b (f(x) - S^*(x))^2 dx \\ &= \int_a^b (S(x) - S^*(x))^2 dx \\ &\quad + 2 \int_a^b (S(x) - S^*(x)) [f(x) - S^*(x)] dx. \end{aligned}$$

由于 $S^*(x)$ 的系数 a_k^* 是方程(4.3)的解, 故

$$\int_a^b (f(x) - S^*(x)) \varphi_k(x) dx = 0 \quad (k = 0, 1, \dots, n),$$

从而上式第二个积分为 0, 于是

$$D = \int_a^b [S(x) - S^*(x)]^2 dx = 0,$$

故(4.4)成立. 这就证明了 $S^*(x)$ 是 $f(x)$ 在 $[a, b]$ 中的最佳平方逼近函数.

若令 $(x) = f(x) - S^*(x)$, 则平方误差为

$$\begin{aligned} \|x\|_2^2 &= (f(x) - S^*(x), f(x) - S^*(x)) \\ &= (f(x), f(x)) - (S^*(x), f(x)) \\ &= \|f(x)\|_2^2 - \sum_{k=0}^n a_k^* (\varphi_k(x), f(x)). \end{aligned} \quad (4.5)$$

若取 $\varphi_k(x) = x^k$, $(x) = 1$, $f(x) \in C[0, 1]$, 则要在 H_n 中求 n 次最佳平方逼近多项式

$$S^*(x) = a_0^* + a_1^* x + \dots + a_n^* x^n,$$

此时

$$\begin{aligned} (\varphi_j(x), \varphi_k(x)) &= \int_0^1 x^{k+j} dx = \frac{1}{k+j+1}, \\ (f(x), \varphi_k(x)) &= \int_0^1 f(x) x^k dx = d_k. \end{aligned}$$

若用 \mathbf{H} 表示 $\mathbf{G}_n = \mathbf{G}(1, x, \dots, x^n)$ 对应的矩阵, 即

$$\mathbf{H} = \begin{matrix} 1 & 1/2 & \dots & 1/(n+1) \\ 1/2 & 1/3 & \dots & 1/(n+2) \\ \dots & \dots & \dots & \dots \\ 1/(n+1) & 1/(n+2) & \dots & 1/(2n+1) \end{matrix} \quad (4.6)$$

称为希尔伯特(Hilbert)矩阵, 记 $\mathbf{a} = (a_0, a_1, \dots, a_n)^T$, $\mathbf{d} = (d_0, d_1, \dots, d_n)^T$, 则

$$\mathbf{H}\mathbf{a} = \mathbf{d} \quad (4.7)$$

的解 $a_k = a_k^*$ ($k=0, 1, \dots, n$) 即为所求.

例 5 设 $f(x) = 1 + x^2$, 求 $[0, 1]$ 上的一次最佳平方逼近多项式.

解 利用(4.7), 得

$$\begin{aligned} d_0 &= \int_0^1 1 + x^2 dx = \frac{1}{2} \ln(1+2) + \frac{2}{2} = 1.147, \\ d_1 &= \int_0^1 x(1+x^2) dx = \frac{1}{3}(1+x^2)^{3/2} \Big|_0^1 = \frac{2\sqrt{2}-1}{3} \\ &\quad 0.609, \end{aligned}$$

得方程组

$$\begin{array}{rcl} 1 & \frac{1}{2} & a \\ & & = \\ \frac{1}{2} & \frac{1}{3} & a \end{array} \quad \begin{array}{l} 1.147 \\ \\ 0.609 \end{array},$$

解出

$$a_0 = 0.934, \quad a_1 = 0.426,$$

故

$$S_1^*(x) = 0.934 + 0.426x.$$

平方误差

$$\begin{aligned} (x) &= (f(x), f(x)) - (S_1^*(x), f(x)) \\ &= \int_0^1 (1+x^2) dx - 0.426 d_1 - 0.934 d_0 = 0.0026. \end{aligned}$$

最大误差

$$(x) = \max_{0 \leq x \leq 1} / 1 + x^2 - S_1^*(x) / 0.066.$$

用 $\{1, x, \dots, x^n\}$ 做基, 求最佳平方逼近多项式, 当 n 较大时, 系数矩阵(4.6)是高度病态的(见第5章), 因此直接求解法方程是相当困难的, 通常是采用正交多项式做基.

3.4.2 用正交函数族作最佳平方逼近

设 $f(x) \in C[a, b]$, $\mathcal{B} = \text{span}\{b_0(x), b_1(x), \dots, b_n(x)\}$, 若 $b_0(x), b_1(x), \dots, b_n(x)$ 是满足条件(2.2)的正交函数族, 则 $(b_i(x), b_j(x)) = 0, i \neq j$ 而 $(b_j(x), b_j(x)) > 0$, 故法方程(4.3)的系数矩阵 $\mathbf{G} = \mathbf{G}(b_0(x), b_1(x), \dots, b_n(x))$ 为非奇异对角阵, 且方程(4.3)的解为

$$a_k^* = (f(x), b_k(x)) / (b_k(x), b_k(x)) \quad (k = 0, 1, \dots, n). \quad (4.8)$$

于是 $f(x) \in C[a, b]$ 在 \mathcal{B} 中的最佳平方逼近函数为

$$S_n^*(x) = \sum_{k=0}^n \frac{(f(x), b_k(x))}{(b_k(x), b_k(x))} b_k(x). \quad (4.9)$$

由(4.5)可得均方误差为

$$\begin{aligned} \|f(x) - S_n^*(x)\|_2^2 &= \|f(x) - S_n^*(x)\|_2^2 \\ &= \|f(x)\|_2^2 - \sum_{k=0}^n \frac{(f(x), b_k(x))^2}{(b_k(x), b_k(x))}. \end{aligned} \quad (4.10)$$

由此可得贝塞尔(Bessel)不等式

$$\sum_{k=1}^n (a_k^* b_k(x))_2^2 \leq \|f(x)\|_2^2. \quad (4.11)$$

若 $f(x) \in C[a, b]$, 按正交函数族 $\{b_k(x)\}$ 展开, 系数 $a_k^* (k=0, 1, \dots)$ 按(4.8)计算, 得级数

$$\sum_{k=0}^{\infty} a_k^* b_k(x) \quad (4.12)$$

称为 $f(x)$ 的广义傅里叶(Fourier)级数, 系数 a_k^* 称为广义傅里叶系数. 它是傅里叶级数的直接推广.

下面讨论特殊情况, 设 $\{b_0(x), b_1(x), \dots, b_n(x)\}$ 是正交多项式, $\mathcal{B} = \text{span}\{b_0(x), b_1(x), \dots, b_n(x)\}, b_k(x) (k=0, 1, \dots, n)$ 可由 $1, x, \dots, x^n$ 正交化得到, 则有下面的收敛定理.

定理7 设 $f(x) \in C[a, b]$, $S_n^*(x)$ 是由(4.9)给出的 $f(x)$ 的最佳平方逼近多项式, 其中 $\{P_k(x), k=0, 1, \dots, n\}$ 是正交多项式族, 则有

$$\lim_{n \rightarrow \infty} \|f(x) - S_n^*(x)\|_2 = 0.$$

证明略, 可见[6].

下面考虑函数 $f(x) \in C[-1, 1]$, 按勒让德多项式 $\{P_0(x), P_1(x), \dots, P_n(x)\}$ 展开, 由(4.8), (4.9)可得

$$S_n^*(x) = a_0^* P_0(x) + a_1^* P_1(x) + \dots + a_n^* P_n(x), \quad (4.13)$$

其中

$$a_k^* = \frac{(f(x), P_k(x))}{(P_k(x), P_k(x))} = \frac{2k+1}{2} \int_{-1}^1 f(x) P_k(x) dx. \quad (4.14)$$

根据(4.10), 平方误差为

$$e_k(x) \|_2^2 = \int_{-1}^1 f^2(x) dx - \sum_{k=0}^n \frac{2}{2k+1} a_k^{*2}. \quad (4.15)$$

由定理7可得

$$\lim_{n \rightarrow \infty} \|f(x) - S_n^*(x)\|_2 = 0.$$

如果 $f(x)$ 满足光滑性条件还可得到 $S_n^*(x)$ 一致收敛于 $f(x)$ 的结论.

定理8 设 $f(x) \in C^2[-1, 1]$, $S_n^*(x)$ 由(4.13)给出, 则对任意 $x \in [-1, 1]$ 和 $\epsilon > 0$, 当 n 充分大时有

$$\|f(x) - S_n^*(x)\| \leq \frac{\epsilon}{n}.$$

证明可见[6].

对于首项系数为1的勒让德多项式 $P_n(x)$ (由公式(2.6)给出)有以下性质.

定理9 在所有最高次项系数为1的 n 次多项式中, 勒让德多项式 $P_n(x)$ 在 $[-1, 1]$ 上与零的平方误差最小.

证明 设 $Q_n(x)$ 是任意一个最高次项系数为 1 的 n 次多项式, 它可表示为

$$Q_n(x) = \varphi_n(x) + \sum_{k=0}^{n-1} a_k \varphi_k(x),$$

于是

$$\begin{aligned} Q_n(x)^2 &= (Q_n(x), Q_n(x)) = \int_{-1}^1 Q_n^2(x) dx \\ &= (\varphi_n(x), \varphi_n(x)) + \sum_{k=0}^{n-1} a_k^2 (\varphi_k(x), \varphi_k(x)) \\ &\quad (\varphi_n(x), \varphi_n(x)) \\ &= \varphi_n(x)^2. \end{aligned}$$

当且仅当 $a_0 = a_1 = \dots = a_{n-1} = 0$ 时等号才成立, 即当 $Q_n(x)$ 是 $\varphi_n(x)$ 时平方误差最小.

例 6 求 $f(x) = e^x$ 在 $[-1, 1]$ 上的三次最佳平方逼近多项式.

解 先计算 $(f(x), \varphi_k(x))$ ($k=0, 1, 2, 3$).

$$(f(x), P_0(x)) = \int_{-1}^1 e^x dx = e - \frac{1}{e} = 2.3504;$$

$$(f(x), P_1(x)) = \int_{-1}^1 x e^x dx = 2e^{-1} = 0.7358;$$

$$(f(x), P_2(x)) = \int_{-1}^1 \frac{3}{2} x^2 - \frac{1}{2} e^x dx = e - \frac{7}{e} = 0.1431;$$

$$(f(x), P_3(x)) = \int_{-1}^1 \frac{5}{2} x^3 - \frac{3}{2} x e^x dx = 37 \frac{1}{e} - 5e = 0.02013.$$

由(4.14)得

$$a^* = (f(x), P_0(x))/2 = 1.1752,$$

$$a^* = 3(f(x), P_1(x))/2 = 1.1036,$$

$$a^* = 5(f(x), P_2(x))/2 = 0.3578,$$

$$a^* = 7(f(x), P_3(x))/2 = 0.07046.$$

代入(4.13)得

$$S_3^*(x) = 0.9963 + 0.9979x + 0.5367x^2 + 0.1761x^3.$$

均方误差

$$\begin{aligned} S_n(x) - S_3^*(x) &= e^x - S_3^*(x) = \int_{-1}^1 e^{2x} dx - \sum_{k=0}^3 \frac{2}{2k+1} a_k^{*2} \\ &= 0.0084. \end{aligned}$$

最大误差

$$S_n(x) = e^x - S_3(x) = 0.0112.$$

如果 $f(x) \in C[a, b]$, 求 $[a, b]$ 上的最佳平方逼近多项式, 做变换

$$x = \frac{b-a}{2}t + \frac{b+a}{2} \quad (-1 \leq t \leq 1),$$

于是 $F(t) = f\left(\frac{b-a}{2}t + \frac{b+a}{2}\right)$ 在 $[-1, 1]$ 上可用勒让德多项式做最佳平方逼近多项式 $S_n^*(t)$, 从而得到区间 $[a, b]$ 上的最佳平方逼近多项式 $S_n^* \frac{1}{b-a}(2x - a - b)$.

由于勒让德多项式 $\{P_k(x)\}$ 是在区间 $[-1, 1]$ 上由 $\{1, x, \dots, x^k, \dots\}$ 正交化得到的, 因此利用函数的勒让德展开部分和得到最佳平方逼近多项式与由

$$S^*(x) = a_0 + a_1 x + \dots + a_n x^n$$

直接通过解法方程得到 H_n 中的最佳平方逼近多项式是一致的, 只是当 n 较大时法方程出现病态, 计算误差较大, 不能使用, 而用勒让德展开不用解线性方程组, 不存在病态问题, 计算公式比较方便, 因此通常都用这种方法求最佳平方逼近多项式.

3.5 曲线拟合的最小二乘法

3.5.1 最小二乘法及其计算

在函数的最佳平方逼近中 $f(x) \in C[a, b]$, 如果 $f(x)$ 只

组离散点集 $\{x_i, i=0, 1, \dots, m\}$ 上给定, 这就是科学实验中经常见到的实验数据 $\{(x_i, y_i), i=0, 1, \dots, m\}$ 的曲线拟合, 这里 $y_i = f(x_i), i=0, 1, \dots, m$, 要求一个函数 $y = S^*(x)$ 与所给数据 $\{(x_i, y_i), i=0, 1, \dots, m\}$ 拟合, 若记误差 $e_i = S^*(x_i) - y_i, i=0, 1, \dots, m$, $e = (e_0, e_1, \dots, e_m)^T$, 设 $\phi_0(x), \phi_1(x), \dots, \phi_n(x)$ 是 $C[a, b]$ 上线性无关函数族, 在 $\mathcal{S} = \text{span}\{\phi_0(x), \phi_1(x), \dots, \phi_n(x)\}$ 中找一函数 $S^*(x)$, 使误差平方和

$$\begin{aligned} \frac{1}{2} \sum_{i=0}^m e_i^2 &= \sum_{i=0}^m [S^*(x_i) - y_i]^2 \\ &= \min_{S(x)} \sum_{i=0}^m [S(x_i) - y_i]^2, \end{aligned} \quad (5.1)$$

这里

$$S(x) = a_0 \phi_0(x) + a_1 \phi_1(x) + \dots + a_n \phi_n(x) \quad (n < m). \quad (5.2)$$

这就是一般的最小二乘逼近, 用几何语言说, 就称为曲线拟合的最小二乘法.

用最小二乘法求拟合曲线时, 首先要确定 $S(x)$ 的形式. 这不单纯是数学问题, 还与所研究问题的运动规律及所得观测数据 (x_i, y_i) 有关; 通常要从问题的运动规律及给定数据描图, 确定 $S(x)$ 的形式, 并通过实际计算选出较好的结果——这点将从下面的例题得到说明. $S(x)$ 的一般表达式为 (5.2) 表示的线性形式. 若 $\psi_k(x)$ 是 k 次多项式, $S(x)$ 就是 n 次多项式. 为了使问题的提法更有一般性, 通常在最小二乘法中 $\frac{1}{2} \sum_{i=0}^m e_i^2$ 都考虑为加权平方和

$$\frac{1}{2} \sum_{i=0}^m w_i(x_i) [S(x_i) - f(x_i)]^2. \quad (5.3)$$

这里 $w(x)$ 是 $[a, b]$ 上的权函数, 它表示不同点 $(x_i, f(x_i))$ 处的数据比重不同, 例如, $w(x_i)$ 可表示在点 $(x_i, f(x_i))$ 处重复观测的次数. 用最小二乘法求拟合曲线的问题, 就是在形如 (5.2) 的

$S(x)$ 中求一函数 $y = S^*(x)$, 使(5.3)取得最小. 它转化为求多元函数

$$\begin{aligned} I(a_0, a_1, \dots, a_n) \\ = \sum_{i=0}^m (x_i) \left[\sum_{j=0}^n a_{j-i} (x_i) - f(x_i) \right]^2 \end{aligned} \quad (5.4)$$

的极小点 $(a_0^*, a_1^*, \dots, a_n^*)$ 问题. 这与第4节讨论的问题完全类似. 由求多元函数极值的必要条件, 有

$$\frac{\partial I}{\partial a_k} = 2 \sum_{i=0}^m (x_i) \left[\sum_{j=0}^n a_{j-i} (x_i) - f(x_i) \right] {}_{-k}(x_i) = 0 \quad (k = 0, 1, \dots, n).$$

若记

$$\begin{aligned} ({}_{-j}, {}_{-k}) &= \sum_{i=0}^m (x_i) {}_{-j}(x_i) {}_{-k}(x_i), \\ (f, {}_{-k}) &= \sum_{i=0}^m (x_i) f(x_i) {}_{-k}(x_i) \quad d_k \\ (k &= 0, 1, \dots, n). \end{aligned} \quad (5.5)$$

上式可改写为

$$\sum_{j=0}^n ({}_{-k}, {}_{-j}) a_j = d_k \quad (k = 0, 1, \dots, n). \quad (5.6)$$

这方程称为法方程, 可写成矩阵形式

$$\mathbf{G}\mathbf{a} = \mathbf{d}.$$

其中 $\mathbf{a} = (a_0, a_1, \dots, a_n)^T$, $\mathbf{d} = (d_0, d_1, \dots, d_n)^T$,

$$\mathbf{G} = \begin{matrix} (0, 0) & (0, 1) & \dots & (0, n) \\ (1, 0) & (1, 1) & \dots & (1, n) \\ \dots & \dots & \dots & \dots \\ (n, 0) & (n, 1) & \dots & (n, n) \end{matrix}. \quad (5.7)$$

要使法方程(5.6)有唯一解 a_0, a_1, \dots, a_n , 就要求矩阵 \mathbf{G} 非奇异, 必须指出, ${}_0(x), {}_1(x), \dots, {}_n(x)$ 在 $[a, b]$ 上线性无关不能推出矩阵 \mathbf{G} 非奇异. 例如, 令 ${}_0(x) = \sin x, {}_1(x) = \sin 2x, x$

$[0, 2]$, 显然 $\{e_0(x), e_1(x)\}$ 在 $[0, 2]$ 上线性无关, 但若取点 $x_k = k$, $k = 0, 1, 2$ ($n = 1, m = 2$), 那么有 $e_0(x_k) = e_1(x_k) = 0$, $k = 0, 1, 2$, 由此得出

$$\mathbf{G} = \begin{pmatrix} (0, 0) & (0, 1) \\ (1, 0) & (1, 1) \end{pmatrix} = 0.$$

为保证(5.6)的系数矩阵 \mathbf{G} 非奇异, 必须加上另外的条件.

定义 10 设 $e_0(x), e_1(x), \dots, e_n(x)$ $C[a, b]$ 的任意线性组合在点集 $\{x_i, i = 0, 1, \dots, m\}$ ($m \geq n$) 上至多只有 n 个不同的零点, 则称 $e_0(x), e_1(x), \dots, e_n(x)$ 在点集 $\{x_i, i = 0, 1, \dots, m\}$ 上满足哈尔(Haar)条件.

显然 $1, x, \dots, x^n$ 在任意 m ($m \geq n$) 个点上满足哈尔条件.

可以证明, 如果 $e_0(x), e_1(x), \dots, e_n(x)$ $C[a, b]$ 在 $\{x_i\}_0^m$ 上满足 Haar 条件, 则方程(5.6)的系数矩阵(5.7)非奇异, 于是方程(5.6)存在唯一的解 $a_k = a_k^*$, $k = 0, 1, \dots, n$. 从而得到函数 $f(x)$ 的最小二乘解为

$$S^*(x) = a_0^* e_0(x) + a_1^* e_1(x) + \dots + a_n^* e_n(x).$$

可以证明这样得到的 $S^*(x)$, 对任何形如(5.2)的 $S(x)$, 都有

$$\sum_{i=0}^m (x_i) [S^*(x_i) - f(x_i)]^2 = \sum_{i=0}^m (x_i) [S(x_i) - f(x_i)]^2,$$

故 $S^*(x)$ 确是所求最小二乘解. 它的证明与(4.4)式相似, 读者可自己完成.

给定 $f(x)$ 的离散数据 $\{(x_i, y_i), i = 0, 1, \dots, m\}$, 要确定 是困难的, 一般可取 $= \text{span}\{1, x, \dots, x^n\}$, 但这样做当 $n \geq 3$ 时, 求解方程(5.6)与连续情形一样, 将出现系数矩阵 \mathbf{G} 为病态的问题, 通常对 $n=1$ 的简单情形都可通过求方程(5.6)得到 $S^*(x)$. 有时根据给定数据图形, 其拟合函数 $y = S(x)$ 表面上不是(5.2)的形式, 但通过变换仍可化为线性模型. 例如, $S(x) = ae^{bx}$, 若两边取对数得

$$\ln S(x) = \ln a + bx,$$

它就是形如(5.2)的线性模型, 具体做法见例8.

例7 已知一组实验数据如下, 求它的拟合曲线.

x_i	1	2	3	4	5
f_i	4	4.5	6	8	8.5
i	2	1	3	1	1

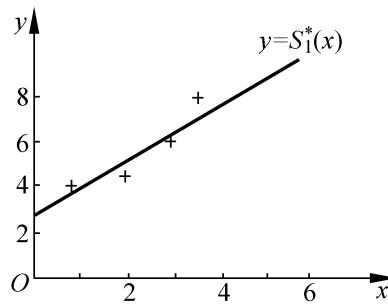


图 3-4

解 根据所给数据, 在坐标纸上标出, 见图3-4. 从图中看到各点在一条直线附近, 故可选择线性函数作拟合曲线, 即令 $S_1(x) = a_0 + a_1 x$, 这里 $m=4$, $n=1$, $a_0(x)=1$, $a_1(x)=x$, 故

$$\begin{aligned} (a_0, a_0) &= \sum_{i=0}^4 a_i = 8, & (a_0, a_1) &= (a_1, a_0) = \sum_{i=0}^4 i x_i = 22, \\ (a_1, a_1) &= \sum_{i=0}^4 i x_i^2 = 74, & (a_0, f) &= \sum_{i=0}^4 i f_i = 47, \\ (a_1, f) &= \sum_{i=0}^4 i x_i f_i = 145.5. \end{aligned}$$

由(5.6)得方程组

$$8a_0 + 22a_1 = 47,$$

$$22a_0 + 74a_1 = 145.5.$$

解得 $a_0 = 2.77$, $a_1 = 1.13$. 于是所求拟合曲线为

$$S_1^*(x) = 2.77 + 1.13x.$$

例8 设数据 (x_i, y_i) ($i = 0, 1, 2, 3, 4$) 由表 3-1 给出, 表中第 4 行为 $\ln y_i = \bar{y}_i$, 可以看出数学模型为 $y = ae^{bx}$, 用最小二乘法确定 a 及 b .

解 根据给定数据 (x_i, y_i) ($i = 0, 1, 2, 3, 4$) 描图可确定拟合曲线方程为 $y = ae^{bx}$, 它不是线性形式. 两边取对数得 $\ln y = \ln a + bx$, 若令 $\bar{y} = \ln y$, $A = \ln a$, 则得 $\bar{y} = A + bx$, $x = \{1, x\}$. 为确定 A , b , 先将 (x_i, y_i) 转化为 (x_i, \bar{y}_i) , 数据表见表 3-1.

表 3-1

i	0	1	2	3	4
x_i	1.00	1.25	1.50	1.75	2.00
y_i	5.10	5.79	6.53	7.45	8.46
\bar{y}_i	1.629	1.756	1.876	2.008	2.135

根据最小二乘法, 取 $\phi_0(x) = 1$, $\phi_1(x) = x$, $\phi_2(x) = 1$, 得

$$(\phi_0, \phi_0) = \sum_{i=0}^4 1 = 5,$$

$$(\phi_0, \phi_1) = \sum_{i=0}^4 x_i = 7.5,$$

$$(\phi_1, \phi_1) = \sum_{i=0}^4 x_i^2 = 11.875,$$

$$(\phi_0, \bar{y}) = \sum_{i=0}^4 \bar{y}_i = 9.404,$$

$$(\phi_1, \bar{y}) = \sum_{i=0}^4 x_i \bar{y}_i = 14.422.$$

故有法方程

$$5A + 7.50b = 9.404,$$

$$7.50A + 11.875b = 14.422.$$

解得 $A = 1.122$, $b = 0.505$, $a = e^A = 3.071$. 于是得最小二乘拟合曲线为

$$y = 3.071e^{0.505x}.$$

现在很多计算机配有自动选择数学模型的程序, 其方法与本例同. 程序中因变量与自变量变换的函数类型较多, 通过计算比较误差找到拟合得较好的曲线, 最后输出曲线图形及数学表达式.

3.5.2 用正交多项式做最小二乘拟合

用最小二乘法得到的法方程组(5.6), 其系数矩阵 \mathbf{G} 是病态的, 但如果 $\phi_0(x), \phi_1(x), \dots, \phi_n(x)$ 是关于点集 $\{x_i\}$ ($i = 0, 1, \dots, m$) 带权 $\omega(x_i)$ ($i = 0, 1, \dots, m$) 正交的函数族, 即

$$\left(\phi_j, \phi_k \right) = \sum_{i=0}^m \phi_j(x_i) \phi_k(x_i) = \begin{cases} 0, & j \neq k, \\ A_k > 0, & j = k, \end{cases} \quad (5.8)$$

则方程(5.6)的解为

$$a_k^* = \frac{\left(f, \phi_k \right)}{\left(\phi_k, \phi_k \right)} = \frac{\sum_{i=0}^m \phi_k(x_i) f(x_i)}{\sum_{i=0}^m \phi_k^2(x_i)} \quad (k = 0, 1, \dots, n), \quad (5.9)$$

且平方误差为

$$f^2 = \sum_{k=0}^n A_k (a_k^*)^2.$$

现在我们根据给定节点 x_0, x_1, \dots, x_m 及权函数 $\omega(x) > 0$, 造出带权 $\omega(x)$ 正交的多项式 $\{P_n(x)\}$. 注意 $n > m$, 用递推公式表示 $P_k(x)$, 即

$$\begin{aligned} P_0(x) &= 1, \\ P_1(x) &= (x - x_0) P_0(x), \\ P_{k+1}(x) &= (x - x_{k+1}) P_k(x) - x_k P_{k-1}(x) \quad (k = 1, 2, \dots, n-1). \end{aligned} \quad (5.10)$$

这里 $P_k(x)$ 是首项系数为 1 的 k 次多项式, 根据 $P_k(x)$ 的正交

性, 得

$$\begin{aligned}
 k+1 &= \frac{\sum_{i=0}^m (x_i) x_i P_k^2(x_i)}{\sum_{i=0}^m (x_i) P_k^2(x_i)} = \frac{(xP_k(x), P_k(x))}{(P_k(x), P_k(x))} \\
 &= \frac{(xP_k, P_k)}{(P_k, P_k)} \quad (k = 0, 1, \dots, n-1), \quad (5.11) \\
 k &= \frac{\sum_{i=0}^m (x_i) P_k^2(x_i)}{\sum_{i=0}^m (x_i) P_{k-1}^2(x_i)} = \frac{(P_k, P_k)}{(P_{k-1}, P_{k-1})} \\
 &\quad (k = 1, \dots, n-1).
 \end{aligned}$$

下面用归纳法证明这样给出的 $\{P_k(x)\}$ 是正交的. 由(5.10)第二式及(5.11)中₁ 的表达式, 有

$$\begin{aligned}
 (P_0, P_1) &= (P_0, xP_0) - (P_0, P_0) \\
 &= (P_0, xP_0) - \frac{(xP_0, P_0)}{(P_0, P_0)} (P_0, P_0) = 0.
 \end{aligned}$$

现假定 $(P_l, P_s) = 0$ ($l < s$) 对 $s = 0, 1, \dots, l-1$ 及 $l = 0, 1, \dots, k$; $k < n$ 均成立, 要证 $(P_{k+1}, P_s) = 0$ 对 $s = 0, 1, \dots, k$ 均成立. 由(5.10)有

$$\begin{aligned}
 (P_{k+1}, P_s) &= ((x - x_{k+1}) P_k, P_s) - (P_{k-1}, P_s) \\
 &= (xP_k, P_s) - (x_{k+1} P_k, P_s) - (P_{k-1}, P_s). \quad (5.12)
 \end{aligned}$$

由归纳法假定, 当 $0 \leq s \leq k-2$ 时

$$(P_k, P_s) = 0, \quad (P_{k-1}, P_s) = 0.$$

另外, $xP_s(x)$ 是首项系数为 1 的 $s+1$ 次多项式, 它可由 P_0, P_1, \dots, P_{s+1} 的线性组合表示, 而 $s+1 \leq k-1$, 故由归纳法假定又有

$$(xP_k, P_s) - (P_k, xP_s) = 0,$$

于是由(5.12), 当 $s \leq k-2$ 时, $(P_{k+1}, P_s) = 0$.

再看

$$\begin{aligned} (P_{k+1}, P_{k-1}) &= (xP_k, P_{k-1}) - {}_{k+1}(P_k, P_{k-1}) \\ &\quad - {}_k(P_{k-1}, P_{k-1}), \end{aligned} \quad (5.13)$$

由假定有

$$(P_k, P_{k-1}) = 0,$$

$$(xP_k, P_{k-1}) = (P_k, xP_{k-1}) = P_k, P_k + \sum_{j=0}^{k-1} c_j P_j = (P_k, P_k).$$

利用(5.11)中_k表达式及以上结果, 得

$$\begin{aligned} (P_{k+1}, P_{k-1}) &= (xP_k, P_{k-1}) - {}_k(P_{k-1}, P_{k-1}) \\ &= (P_k, P_k) - (P_k, P_k) = 0. \end{aligned}$$

最后, 由(5.11)有

$$\begin{aligned} (P_{k+1}, P_k) &= (xP_k, P_k) - {}_{k+1}(P_k, P_k) - {}_k(P_k, P_{k-1}) \\ &= (xP_k, P_k) - \frac{(xP_k, P_k)}{(P_k, P_k)} (P_k, P_k) = 0. \end{aligned}$$

至此已证明了由(5.10)及(5.11)确定的多项式 $\{P_k(x)\}$ ($k=0, 1, \dots, n, n-m$)组成一个关于点集 $\{x_i\}$ 的正交系.

用正交多项式 $\{P_k(x)\}$ 的线性组合作最小二乘曲线拟合, 只要根据公式(5.10)及(5.11)逐步求 $P_k(x)$ 的同时, 相应计算出系数

$$a_k = \frac{(f, P_k)}{(P_k, P_k)} = \frac{\sum_{i=0}^m (x_i) f(x_i) P_k(x_i)}{\sum_{i=0}^m (x_i) P_k^2(x_i)} \quad (k=0, 1, \dots, n),$$

并逐步把 $a_k^* P_k(x)$ 累加到 $S(x)$ 中去, 最后就可得到所求的拟合曲线

$$y = S(x) = a_0^* P_0(x) + a_1^* P_1(x) + \dots + a_n^* P_n(x).$$

这里 n 可事先给定或在计算过程中根据误差确定. 用这种方法编程序不用解方程组, 只用递推公式, 并且当逼近次数增加一次时, 只要把程序中循环数加1, 其余不用改变. 这是目前用多项式做曲线拟合最好的计算方法, 有通用的语言程序供用户使用.

3.6 最佳平方三角逼近与快速傅里叶变换

当 $f(x)$ 是周期函数时, 显然用三角多项式逼近 $f(x)$ 比用代数多项式更合适, 本节主要讨论用三角多项式做最小平方逼近及快速傅里叶变换(Fast Fourier Transform)简称 FFT 算法.

3.6.1 最佳平方三角逼近与三角插值

设 $f(x)$ 是以 2π 为周期的平方可积函数, 用三角多项式

$$S_n(x) = \frac{1}{2}a_0 + a_1 \cos x + b_1 \sin x + \dots + a_n \cos nx + b_n \sin nx \quad (6.1)$$

做最佳平方逼近函数. 由于三角函数族

$$1, \cos x, \sin x, \dots, \cos kx, \sin kx, \dots$$

在 $[0, 2\pi]$ 上是正交函数族, 于是 $f(x)$ 在 $[0, 2\pi]$ 上的最小平方三角逼近多项式 $S_n(x)$ 的系数是

$$\begin{aligned} a_k &= \frac{1}{2} \int_0^{2\pi} f(x) \cos kx dx \quad (k = 0, 1, \dots, n), \\ b_k &= \frac{1}{2} \int_0^{2\pi} f(x) \sin kx dx \quad (k = 1, \dots, n), \end{aligned} \quad (6.2)$$

a_k, b_k 称为傅里叶系数, 函数 $f(x)$ 按傅里叶系数展开得到的级数

$$\frac{1}{2}a_0 + \sum_{k=1}^n (a_k \cos kx + b_k \sin kx) \quad (6.3)$$

就称为傅里叶级数, 只要 $f(x)$ 在 $[0, 2\pi]$ 上分段连续, 则级数(6.3)一致收敛到 $f(x)$.

对于最佳平方逼近多项式(6.1)有

$$\|f(x) - S_n(x)\|_2^2 = \|f(x) - S_n(x)\|_2^2 - \|S_n(x)\|_2^2.$$

由此可以得到相应于(4.11)的贝塞尔不等式

$$\frac{1}{2} a^2 + \sum_{k=1}^n (a_k^2 + b_k^2) - \frac{1}{0} [f(x)]^2 dx.$$

因为右边不依赖于 n , 左边单调有界, 所以级数 $\frac{1}{2} a^2 + \sum_{k=1}^n (a_k^2 + b_k^2)$ 收敛, 并有

$$\lim_k a_k = \lim_k b_k = 0.$$

当 $f(x)$ 只在给定的离散点集 $x_j = \frac{2}{N}j, j = 0, 1, \dots, N-1$ 上

已知时, 则可类似得到离散点集正交性与相应的离散傅里叶系数. 为方便起见, 下面只给出奇数个点的情形. 令

$$x_j = \frac{2j}{2m+1} \quad (j = 0, 1, \dots, 2m),$$

可以证明对任何 $0 < k, l < m$ 成立

$$\sum_{j=0}^{2m} \sin lx_j \sin kx_j = 0, \quad l = k, l = k = 0;$$

$$\sum_{j=0}^{2m} \cos lx_j \cos kx_j = \frac{2m+1}{2}, \quad l = k = 0;$$

$$0, \quad l = k;$$

$$\sum_{j=0}^{2m} \cos lx_j \cos kx_j = \frac{2m+1}{2}, \quad l = k = 0;$$

$$2m+1, \quad l = k = 0;$$

$$\sum_{j=0}^{2m} \cos lx_j \sin kx_j = 0, \quad 0 < k, j < m.$$

这就表明函数族 $\{1, \cos x, \sin x, \dots, \cos mx, \sin mx\}$ 在点集

$x_j = \frac{2j}{2m+1}$ 上正交, 若令 $f_j = f(x_j)$ ($j = 0, 1, \dots, 2m$), 则

$f(x)$ 的最小二乘三角逼近为

$$S_n(x) = \frac{1}{2} a_0 + \sum_{k=1}^n (a_k \cos kx + b_k \sin kx), \quad n < m,$$

其中

$$\begin{aligned} a_k &= \frac{2}{2m+1} \sum_{j=0}^{2m} f_j \cos \frac{2jk}{2m+1} \quad (k = 0, 1, \dots, m), \\ b_k &= \frac{2}{2m+1} \sum_{j=0}^{2m} f_j \sin \frac{2jk}{2m+1} \quad (k = 1, \dots, n). \end{aligned} \quad (6.4)$$

当 $n = m$ 时, 可证明

$$S_m(x_j) = f_j \quad (j = 0, 1, \dots, 2m),$$

于是

$$S_m(x) = \frac{1}{2} a_0 + \sum_{k=1}^m (a_k \cos kx + b_k \sin kx)$$

就是三角插值多项式, 系数仍由(6.4)表示.

更一般情形, 假定 $f(x)$ 是以 2 为周期的复函数, 给定在 N 个等分点 $x_j = \frac{2}{N}j$ ($j = 0, 1, \dots, N-1$) 上的值 $f_j = f(\frac{2}{N}j)$, 由于 $e^{ijx} = \cos(jx) + i\sin(jx)$ ($j = 0, 1, \dots, N-1$, $i = -1$), 函数族 $\{1, e^{ix}, \dots, e^{i(N-1)x}\}$ 在区间 $[0, 2]$ 上是正交的, 函数 e^{ijx} 在等距点集 $x_k = \frac{2}{N}k$ ($k = 0, 1, \dots, N-1$) 上的值 e^{ijx_k} 组成的向量记作

$$v_j = (1, e^{ij\frac{2}{N}}, \dots, e^{ij\frac{2}{N}(N-1)})^T.$$

当 $j = 0, 1, \dots, N-1$ 时, N 个复向量 v_0, v_1, \dots, v_{N-1} 具有下面所定义的正交性:

$$\begin{aligned} (v_l, v_s) &= \sum_{k=0}^{N-1} e^{il\frac{2}{N}k} e^{-is\frac{2}{N}k} = \sum_{k=0}^{N-1} e^{i(l-s)\frac{2}{N}k} \\ &= \begin{cases} 0, & l = s; \\ N, & l \neq s. \end{cases} \end{aligned} \quad (6.5)$$

事实上, 令 $r = e^{i(l-s)\frac{2}{N}}$, 若 $0 \leq l, s \leq N-1$, 则有

$$0 \leq l \leq N-1, \quad - (N-1) \leq l-s \leq N-1, \quad -s \leq 0,$$

于是 $- (N-1) \leq l-s \leq N-1$,

即 $-1 < -\frac{N-1}{N} \leq \frac{l-s}{N} \leq \frac{N-1}{N} < 1$;

若 $l = s = 0$, 则 $r = 1$, 从而

$$r^N = e^{i(l-s)2\pi} = 1;$$

于是

$$(l, s) = \sum_{k=0}^{N-1} r^k = \frac{1 - r^N}{1 - r} = 0.$$

若 $l = s$, 则 $r = 1$, 于是

$$(s, s) = \sum_{k=0}^{N-1} r^k = N.$$

这就证明了(6.5)成立. 即 e_0, e_1, \dots, e_{N-1} 是正交的.

因此, $f(x)$ 在 N 个点 $x_j = \frac{2\pi}{N}j, j = 0, 1, \dots, N-1$ 上的最小

二乘傅里叶逼近为

$$S(x) = \sum_{k=0}^{n-1} c_k e^{ikx}, \quad n \leq N, \quad (6.6)$$

其中

$$\alpha_k = \frac{1}{N} \sum_{j=0}^{N-1} f_j e^{-ikj\frac{2\pi}{N}} \quad (k = 0, 1, \dots, n-1). \quad (6.7)$$

在(6.6)中若 $n = N$, 则 $S(x)$ 为 $f(x)$ 在点 $x_j (j = 0, 1, \dots, N-1)$ 上的插值函数, 即 $S(x_j) = f(x_j)$, 于是由(6.6)得

$$f_j = \sum_{k=0}^{N-1} c_k e^{ikj\frac{2\pi}{N}} \quad (j = 0, 1, \dots, N-1). \quad (6.8)$$

(6.7)是由 $\{f_j\}$ 求 $\{c_k\}$ 的过程, 称为 $f(x)$ 的离散傅里叶变换. 简称 DFT, 而(6.8)是由 $\{\alpha_k\}$ 求 $\{f_j\}$ 的过程, 称为反变换. 它们是使用计算机进行傅里叶分析(简称傅氏分析)的主要方法, 在数字讯号处理, 全息技术, 光谱和声谱分析等很多领域都有广泛应用.

3.6.2 快速傅氏变换(FFT)

不论是按(6.7)式由 $\{f_k\}$ 求 $\{c_k\}$ 或是按(6.8)式由 $\{c_k\}$ 求 $\{f_k\}$, 还是由(6.4)计算傅里叶逼近系数 a_k, b_k 都可归结为计算

$$c_j = \sum_{k=0}^{N-1} x_k w^{kj} \quad (j = 0, 1, \dots, N-1), \quad (6.9)$$

其中 $w = \exp(-i2\pi/N)$ (正变换) 或 $w = \exp(i2\pi/N)$ (反变换), $\{x_k\}$ ($k = 0, 1, \dots, N-1$) 是已知复数序列. 如直接用(6.9)计算 c_j , 需要 N 次复数乘法和 N 次复数加法, 称为 N 个操作, 计算全部 c_j 共要 N^2 个操作. 当 N 较大且处理数据很多时, 就是用高速的电子计算机, 很多实际问题仍然无法计算, 直到 20 世纪 60 年代中期产生了 FFT 算法, 大大提高了运算速度, 才使傅氏变换得以广泛应用. FFT 算法的思想就是尽量减少乘法次数, 例如计算 $ab + ac = a(b + c)$, 用左端计算时做两次乘法, 右端只用一次乘法. 用公式(6.9)计算全部 c_j , 表面看要做 N^2 个乘法, 实际上所有 $\exp(i2\pi kj/N)$, $j, k = 0, 1, \dots, N-1$ 中, 只有 N 个不同的值 w^0, w^1, \dots, w^{N-1} , 特别当 $N = 2^p$ 时, 只有 $N/2$ 个不同的值. 因此, 我们可把同一个 w^r 对应的 x_k 相加后再乘 w^r 这就能大量减少乘法次数. 为了具体推导 FFT 算法, 先给出定义.

设正整数 m 除以 N 后得商 q 及余数 r , 则 $m = qN + r$, r 称为 m 的 N 同余数, 以 $m \mod N$ 表示. 由于 $w = \exp(i2\pi/N)$, $w^N = e^{i2\pi} = 1$, 故有 $w^m = (w^N)^q w^r = w^r$.

因此计算 w^m 时可用 w 的 N 同余数 r 代替 m , 从而推出 FFT 算法. 下面以 $N = 2^3$ 为例. 说明 FFT 的计算方法. 由于 $0 \leq k, j \leq N-1 = 2^3 - 1 = 7$, 则(6.9)的和是

$$c_j = \sum_{k=0}^7 x_k w^{jk} \quad (j = 0, 1, \dots, 7). \quad (6.10)$$

将 k, j 用二进制表示为

$$\begin{aligned} k &= k_2 2^2 + k_1 2^1 + k_0 2^0 = (k_2 k_1 k_0), \\ j &= j_2 2^2 + j_1 2^1 + j_0 2^0 = (j_2 j_1 j_0); \end{aligned}$$

其中 k_r, j_r ($r = 0, 1, 2$) 只能取 0 或 1, 例如 $6 = 2^2 + 2^1 + 0 \cdot 2^0 = (110)_2$. 根据 k, j 表示法, 有

$$c_j = c(j_2 j_1 j_0), \quad x_k = x(k_2 k_1 k_0).$$

公式(6.10)可表示为

$$\begin{aligned} & c(j_2 j_1 j_0) \\ &= \sum_{\substack{k_2=0 \\ k_1=0 \\ k_0=0}}^{1 \quad 1 \quad 1} x(k_2 k_1 k_0) w^{(k_2 k_1 k_0)(j_2^2 + j_1^2 + j_0^2)} \\ &= \sum_{\substack{k_2=0 \\ k_1=0 \\ k_0=0}}^{1 \quad 1 \quad 1} x(k_2 k_1 k_0) w^{j_0(k_2 k_1 k_0)} \cdot w^{j_1(k_1 k_0)} \\ & \quad \cdot w^{j_2(k_0)} . \end{aligned} \quad (6.11)$$

若引入记号

$$\begin{aligned} A_0(k_2 k_1 k_0) &= x(k_2 k_1 k_0), \\ A_1(k_1 k_0 j_0) &= \sum_{k_2=0}^1 A_0(k_2 k_1 k_0) w^{j_0(k_2 k_1 k_0)}, \\ A_2(k_0 j_1 j_0) &= \sum_{k_1=0}^1 A_1(k_1 k_0 j_0) w^{j_1(k_1 k_0)}, \\ A_3(j_2 j_1 j_0) &= \sum_{k_0=0}^1 A_2(k_0 j_1 j_0) w^{j_2(k_0)}. \end{aligned} \quad (6.12)$$

则(6.11)变成

$$c(j_2 j_1 j_0) = A_3(j_2 j_1 j_0).$$

它说明利用 N 同余数可把计算 c_j 分为 p 步, 用公式(6.12)计算, 每计算一个 A_q 只用 2 次复数乘法, 计算一个 c_j 用 $2p$ 次复数乘法, 计算全部 c_j 共用 $2pN$ 次复数乘法. 若注意 $w^{j_0 2^{p-1}} = w^{j_0 N^2} = (-1)^j$, 公式(6.12)还可进一步简化为

$$\begin{aligned} A_1(k_1 k_0 j_0) &= \sum_{k_2=0}^1 A_0(k_2 k_1 k_0) w^{j_0(k_2 k_1 k_0)} \\ &= A_0(0 k_1 k_0) w^{j_0(0 k_1 k_0)} \\ &\quad + A_0(1 k_1 k_0) w^{j_0 2^2} w^{j_0(0 k_1 k_0)} \\ &= [A_0(0 k_1 k_0) + (-1)^{j_0} A_0(1 k_1 k_0)] w^{j_0(0 k_1 k_0)}, \end{aligned}$$

$$A_1(k_1 k_0) = A_0(0 k_1 k_0) + A_0(1 k_1 k_0),$$

$$A_1(k_1 k_0 1) = [A_0(0 k_1 k_0) - A_0(1 k_1 k_0)] w^{(0 k_1 k_0)}.$$

将这表达式中二进制表示还原为十进制表示: $k = (0 k_1 k_0) = k_1 2^1 + k_0 2^0$, 即 $k = 0, 1, 2, 3$, 得

$$\begin{aligned} A_1(2k) &= A_0(k) + A_0(k + 2^2), \\ A_1(2k + 1) &= [A_0(k) - A_0(k + 2^2)] w^k \quad (6.13) \\ (k &= 0, 1, 2, 3). \end{aligned}$$

同样(6.12)中的 A_2 也可简化为

$$A_2(k_0 j_1 j_0) = [A_1(0 k_0 j_0) + (-1)^{j_1} A_1(1 k_0 j_0)] w^{j_1(0 k_0 0)},$$

即

$$\begin{aligned} A_2(k_0 0 j_0) &= A_1(0 k_0 j_0) + A_1(1 k_0 j_0), \\ A_2(k_0 1 j_0) &= [A_1(0 k_0 j_0) - A_1(1 k_0 j_0)] w^{(0 k_0 0)}. \end{aligned}$$

把二进制表示还原为十进制表示, 得

$$\begin{aligned} A_2(k2^2 + j) &= A_1(2k + j) + A_1(2k + j + 2^2), \\ A_2(k2^2 + j + 2) &= [A_1(2k + j) - A_1(2k + j + 2^2)] w^{2k} \\ (k &= 0, 1; j = 0, 1). \quad (6.14) \end{aligned}$$

同理(6.12)中 A_3 可简化为

$$A_3(j_2 j_1 j_0) = A_2(0 j_1 j_0) + (-1)^{j_2} A_2(1 j_1 j_0),$$

即

$$\begin{aligned} A_3(0 j_1 j_0) &= A_2(0 j_1 j_0) + A_2(1 j_1 j_0), \\ A_3(1 j_1 j_0) &= A_2(0 j_1 j_0) - A_2(1 j_1 j_0). \end{aligned}$$

表示为十进制, 有

$$\begin{aligned} A_3(j) &= A_2(j) + A_2(j + 2^2), \\ A_3(j + 2^2) &= A_2(j) - A_2(j + 2^2) \\ (j &= 0, 1, 2, 3). \quad (6.15) \end{aligned}$$

根据公式(6.13), (6.14), (6.15), 由 $A_0(k) = x(k) = x_k (k = 0, 1, \dots, 7)$ 逐次计算到 $A_3(j) = c_j (j = 0, 1, \dots, 7)$, 见表 3-2.

上面推导的 $N = 2^3$ 的计算公式可类似地推广到 $N = 2^p$ 的情形. 根据公式(6.13), (6.14), (6.15), 一般情况的 FFT 计算公式如下:

$$\begin{aligned} A_q(k2^q + j) &= A_{q-1}(k2^{q-1} + j) + A_{q-1}(k2^{q-1} + j + 2^{p-1}), \\ A_q(k2^q + j + 2^{q-1}) &= [A_{q-1}(k2^{q-1} + j) - A_{q-1}(k2^{q-1} + j + 2^{p-1})] w^{k2^{q-1}}, \end{aligned} \quad (6.16)$$

其中 $q = 1, \dots, p$; $k = 0, 1, \dots, 2^{p-q} - 1$; $j = 0, 1, \dots, 2^{q-1} - 1$. A_q 括号内的数代表它的位置, 在计算机中代表存放数的地址. 一组 A_q 占用 N 个复数单元, 计算时需给出两组单元, 从 $A_0(m)$ ($m = 0, 1, \dots, N-1$) 出发, q 由 1 算到 p , $A_p(j) = c_j$ ($j = 0, 1, \dots, N-1$), 即为所求. 计算过程中只要按地址号存放 A_q , 则最后得到的 $A_p(j)$ 就是所求离散频谱的次序. (注意, 目前一些计算机程序计算结果地址是逆序排列, 还要增加倒地址的一步才是我们这里介绍的结果). 这个计算公式除了具有不倒地址的优点外, 计算只有两重循环, 外循环 q 由 1 计算到 p , 内循环 k 由 0 计算到 $2^{p-q} - 1$, j 由 0 计算到 $2^{q-1} - 1$, 更重要的是整个计算过程省计算量. 由公式看到算一个 A_q 共做 $2^{p-q} 2^{q-1} = N/2$ 次复数乘法, 而最后一步计算 A_p 时, 由于 $w^{k2^{p-1}} = (w^{N/2})^k = (-1)^k = (-1)^0 = 1$ (注意 $q = p$ 时 $2^{p-q} - 1 = 0$, 故 $k = 0$), 因此, 总共要算 $(p-1)N/2$ 次复数乘法, 它比直接用(6.9)需 N^2 次乘法快得多, 计算量比值是 $N:(p-1)/2$. 当 $N = 2^{10}$ 时比值是 $1024/4.5 = 230$, 它比一般 FFT 的计算量 (pN 次乘法) 也快一倍. 我们称(6.16)的计算公式为改进的 FFT 算法, 下面给出这一算法的程序步骤:

步骤 1 给出数组 $A_1(N), A_2(N)$ 及 $w(N/2)$;

步骤 2 将已知的记录复数数组 $\{x_k\}$ 输入到单元 $A_1(k)$ 中 (k 从 0 到 $N-1$);

表 3-2

单元码号	0 000	1 001	2 010	3 011	4 100	5 101	6 110	7 111
					$w^0 = 1$	w^1	w^2	w^3
$x_k = A_0(k)$	$A_0(0)$	$A_0(1)$	$A_0(2)$	$A_0(3)$	$A_0(4)$	$A_0(5)$	$A_0(6)$	$A_0(7)$
A_1	$A_0(0) + [A_0(0) - A_0(4)] w^0$	$A_0(1) + [A_0(1) - A_0(5)] w^1$	$A_0(2) + [A_0(2) - A_0(6)] w^2$	$A_0(3) + [A_0(3) - A_0(7)] w^3$				
A_2	$A_1(0) + [A_1(0) - A_1(4)] w^0$	$A_1(1) + [A_1(1) - A_1(5)] w^1$	$A_1(2) + [A_1(2) - A_1(6)] w^2$	$A_1(3) + [A_1(3) - A_1(7)] w^3$				
$c_i = A_3(j)$	$(0) + (4)$	$(1) + (5)$	$(2) + (6)$	$(3) + (7)$	$(0) - (4)$	$(1) - (5)$	$(2) - (6)$	$(3) - (7)$

步骤 3 计算 $w^m = \exp -i \frac{2}{N} m$ 或 $w^m = \exp i \frac{2}{N} m$ 存

放在单元 $w(m)$ 中 (m 从 0 到 $(N/2) - 1$);

步骤 4 q 循环从 1 到 p , 若 q 为奇数做步骤 5, 否则做步骤 6;

步骤 5 k 循环从 0 到 $2^{p-q} - 1$, j 循环从 0 到 $2^{q-1} - 1$, 计算

$$\begin{aligned} A_2(k2^q + j) &= A_1(k2^{q-1} + j) + A_1(k2^{q-1} + j + 2^{p-1}), \\ A_2(k2^q + j + 2^{q-1}) &= [A_1(k2^{q-1} + j) \\ &\quad - A_1(k2^{q-1} + j + 2^{q-1})] w(k2^{q-1}); \end{aligned}$$

转步骤 7.

步骤 6 k 循环从 0 到 $2^{p-q} - 1$, j 循环从 0 到 2^{q-1} , 计算

$$\begin{aligned} A_1(k2^q + j) &= A_2(k2^{q-1} + j) + A_2(k2^{q-1} + j + 2^{p-1}), \\ A_1(k2^q + j + 2^{q-1}) &= [A_2(k2^{q-1} + j) \\ &\quad - A_2(k2^{q-1} + j + 2^{q-1})] w(k2^{q-1}). \end{aligned}$$

k, j 循环结束, 做下一步;

步骤 7 若 $q = p$ 转步骤 8, 否则 $q+1 = q$ 转步骤 4.

步骤 8 q 循环结束, 若 $p = \text{偶数}$, 将 $A_1(j) - A_2(j)$, 则 $c_j = A_2(j)$ ($j = 0, 1, \dots, N-1$) 即为所求.

3.7 有理逼近

3.7.1 有理逼近与连分式

前面讨论了用多项式逼近函数 $f(x) \in C[a, b]$, 多项式是一种计算简便的函数类, 但当函数在某点附近无界时用多项式逼近效果很差, 而用有理函数逼近则可得到较好的效果. 所谓有理函数逼近是指用形如

$$R_{nm}(x) = \frac{P_n(x)}{Q_m(x)} = \frac{\sum_{k=0}^n a_k x^k}{\sum_{k=0}^m b_k x^k} \quad (7.1)$$

的函数逼近 $f(x)$, 与前面讨论一样, 如果 $\|f(x) - R_{nm}(x)\|_2$ 最小就可得到最佳有理一致逼近, 如果 $\|f(x) - R_{nm}(x)\|_\infty$ 最小则可得到最佳有理平方逼近函数. 这里不做具体介绍, 可参看文献 [5]. 本节主要讨论利用函数的泰勒展开获得有理逼近函数的方法. 先看例题, 对函数 $\ln(1+x)$ 用泰勒展开得

$$\ln(1+x) = \sum_{k=1}^{\infty} (-1)^{k-1} \frac{x^k}{k} \quad x \in [-1, 1]. \quad (7.2)$$

取部分和

$$S_n(x) = \sum_{k=1}^n (-1)^{k-1} \frac{x^k}{k} \quad \ln(1+x).$$

另方面若对(7.2)式用辗转相除可得到 $\ln(1+x)$ 的一种连分式展开

$$\begin{aligned} \ln(1+x) &= \cfrac{x}{1 + \cfrac{1-x}{2 + \cfrac{1-x}{3 + \cfrac{2^2 \cdot x}{4 + \cfrac{2^2 \cdot x}{5 + \dots}}}}} \end{aligned}$$

$$= \frac{x}{1} + \frac{1+x}{2} + \frac{1+x}{3} + \frac{2^2+x}{4} + \frac{2^2+x}{5} + \dots . \quad (7.3)$$

(7.3) 右端为 $\ln(1+x)$ 的无穷连分式的前 5 项, 最后式子是它的紧凑形式, 若取(7.3)的前 2, 4, 6, 8 项, 则可分别得到 $\ln(1+x)$ 的以下有理逼近

$$\begin{aligned} R_{11}(x) &= \frac{2x}{2+x}, & R_{22}(x) &= \frac{6x+3x^2}{6+6x+x^2}, \\ R_{33}(x) &= \frac{60x+60x^2+11x^3}{60+90x+36x^2+3x^3}, & (7.4) \\ R_{44}(x) &= \frac{420x+630x^2+260x^3+25x^4}{420+840x+540x^2+120x^3+6x^4}. \end{aligned}$$

若用同样多项的泰勒展开部分和 $S_{2n}(x)$ 逼近 $\ln(1+x)$, 并计算 $x=1$ 处的值 $S_{2n}(1)$ 及 $R_{nn}(1)$, 计算结果见表 3-3.

表 3-3

n	$S_{2n}(1)$	$s = \ln 2 - S_{2n}(1) $	$R_m(1)$	$r = \ln 2 - R_m(1) $
1	0.50	0.19	0.667	0.026
2	0.58	0.11	0.69231	0.00084
3	0.617	0.076	0.693122	0.000025
4	0.634	0.058	0.69314642	0.00000076

$\ln 2$ 的准确值为 $0.69314718\dots$, 由此看出 $R_{44}(1)$ 的精度比 $S_8(1)$ 高出近 10 万倍, 而它们的计算量是相当的, 这说明用有理逼近比多项式逼近好得多. 在计算机上计算有理函数(7.1)的值通常可转化为连分式, 这样可以节省乘除法的计算次数.

例 9 给出有理函数

$$R_{43}(x) = \frac{2x^4 + 45x^3 + 381x^2 + 1353x + 1511}{x^3 + 21x^2 + 157x + 409}$$

用辗转相除法将它化为连分式并写成紧凑形式.

解 用辗转相除可逐步得到

$$\begin{aligned}
 R_{43}(x) &= 2x + 3 + \frac{4x^2 + 64x + 284}{x^3 + 21x^2 + 157x + 409} \\
 &= 2x + 3 + \frac{4}{x + 5 + \frac{6(x+9)}{x^2 + 16x + 71}} \\
 &= 2x + 3 + \frac{4}{x + 5 + \frac{6}{x + 7 + \frac{8}{x + 9}}} \\
 &= 2x + 3 + \frac{4}{x + 5} + \frac{6}{x + 7} + \frac{8}{x + 9}.
 \end{aligned}$$

本例中用连分式计算 $R_{43}(x)$ 的值只需 3 次除法, 1 次乘法和 7 次加法. 若直接用多项式计算的秦九韶算法则需 6 次乘法和 1 次除法及 7 次加法, 可见将 $R_{nm}(x)$ 化成连分式可节省计算乘除法次数, 对一般的有理函数(7.1)可转化为一个连分式

$$R_{nm}(x) = P_1(x) + \frac{c_2}{x + d_1} + \dots + \frac{c_l}{x + d_l}.$$

它的乘除法运算只需 $\max(m, n)$ 次, 而直接用有理函数(7.1)计算乘除法次数为 $n+m$ 次.

3.7.2 帕德逼近

利用函数 $f(x)$ 的泰勒展开可以得到它的有理逼近. 设 $f(x)$ 在 $x=0$ 的泰勒展开为

$$f(x) = \sum_{k=0}^N \frac{1}{k!} f^{(k)}(0) x^k + \frac{f^{(N+1)}(\cdot)}{(N+1)!} x^{N+1}. \quad (7.5)$$

它的部分和记作

$$P(x) = \sum_{k=0}^N \frac{1}{k!} f^{(k)}(0) x^k = \sum_{k=0}^N c_k x^k. \quad (7.6)$$

定义 11 设 $f \in C^{N+1}(-a, a)$, $N = n + m$, 如果有理函数

$$R_{nm}(x) = \frac{a_0 + a_1 x + \dots + a_n x^n}{1 + b_1 x + \dots + b_m x^m} = \frac{P_n(x)}{Q_m(x)}, \quad (7.7)$$

其中 $P_n(x)$, $Q_m(x)$ 无公因式, 且满足条件

$$R_{nm}^{(k)}(0) = f^{(k)}(0) \quad (k = 0, 1, \dots, N), \quad (7.8)$$

则称 $R_{nm}(x)$ 为函数 $f(x)$ 在 $x=0$ 处的 (n, m) 阶帕德(**Padé**)逼近, 记作 $R(n, m)$, 简称 $R(n, m)$ 的帕德逼近.

根据定义, 若令

$$h(x) = P(x) Q_m(x) - P_n(x),$$

则满足条件(7.8)等价于

$$h^{(k)}(0) = 0, \quad k = 0, 1, \dots, N.$$

即

$$h^{(k)}(0) = (P(x) Q_m(x) - P_n(x))^{(k)} \Big|_{x=0} = 0, \quad k = 0, 1, \dots, N.$$

由于 $P_n^{(k)}(0) = k! a_k$, 应用莱布尼兹求导公式得

$$(P(x) Q_m(x) - P_n(x))^{(k)} \Big|_{x=0} = k! \sum_{j=0}^k c_j b_{k-j} - k! a_k = 0 \\ k = 0, 1, \dots, N,$$

这里 $c_j = \frac{1}{j!} f^{(j)}(0)$ 是由(7.6)得到的, 上式两端除 $k!$ 并由 $b_0 = 1$,

$b_j = 0$ (当 $j > m$ 时), 可得

$$a_k = \sum_{j=0}^{k-1} c_j b_{k-j} + c_k \quad (k = 0, 1, \dots, n) \quad (7.9)$$

及

$$- \sum_{j=0}^{k-1} c_j b_{k-j} = c_k \quad (k = n+1, \dots, n+m). \quad (7.10)$$

注意当 $j > m$ 时 $b_j = 0$, 故(7.10)可写成

$$\begin{aligned} -c_{n-m+1} b_m - \dots - c_{n-1} b_2 - c_n b_1 &= c_{n+1}, \\ -c_{n-m+2} b_m - \dots - c_n b_2 - c_{n+1} b_1 &= c_{n+2}, \\ \dots \\ -c_n b_m - \dots - c_{n+m-2} b_2 - c_{n+m-1} b_1 &= c_{n+m}. \end{aligned} \quad (7.11)$$

其中 $j < 0$ 时 $c_j = 0$, 若记

$$\mathbf{H} = \begin{bmatrix} - & C_{n-m+1} & \cdots & - & C_{n-1} & - & C_n \\ - & C_{n-m+2} & \cdots & - & C_n & - & C_{n+1} \\ \cdots & & & & \cdots & & \cdots \\ - & C_1 & \cdots & - & C_{n+m-2} & - & C_{n+m-1} \end{bmatrix}, \quad (7.12)$$

$$\mathbf{B} = (b_m, b_{m-1}, \dots, b)^T, \quad \mathbf{C} = (c_{n+1}, c_{n+2}, \dots, c_{n+m})^T.$$

则方程组(7.11)的矩阵形式为

$$\mathbf{H}\mathbf{B} = \mathbf{C}.$$

综上所述得下面的定理 .

定理 10 设 $f(x) \in C^{N+1}(-a, a)$, $N = n + m$, 则形如(7.7)的有理函数 $R_{nm}(x)$ 是 $f(x)$ 的 (n, m) 阶帕德逼近的充分必要条件是多项式 $P_n(x)$ 及 $Q_m(x)$ 的系数 a_0, a_1, \dots, a_n 及 b_1, \dots, b_m 满足方程组(7.9)及(7.11) .

根据定理 10 求 $f(x)$ 的帕德逼近时, 首先要由(7.11)解出 $Q_m(x)$ 的系数 b_1, \dots, b_m , 再由(7.9)直接算出 $P_n(x)$ 的系数 a_0, a_1, \dots, a_n . $f(x)$ 的各阶帕德逼近可列成一张表, 称为帕德表(见表 3-4) .

表 3-4 帕德表

$n \backslash m$	0	1	2	3	4	...
0	(0, 0)	(0, 1)	(0, 2)	(0, 3)	(0, 4)	...
1	(1, 0)	(1, 1)	(1, 2)	(1, 3)	(1, 4)	...
2	(2, 0)	(2, 1)	(2, 2)	(2, 3)	(2, 4)	...
3	(3, 0)	(3, 1)	(3, 2)	(3, 3)	(3, 4)	...
4	(4, 0)	(4, 1)	(4, 2)	(4, 3)	(4, 4)	...
...	

例 10 求 $f(x) = \ln(1 + x)$ 的帕德逼近 $R(2, 2)$ 及 $R(3, 3)$.

解 由 $\ln(1+x)$ 的泰勒展开

$$\ln(1+x) = x - \frac{1}{2}x^2 + \frac{1}{3}x^3 - \frac{1}{4}x^4 + \dots$$

得 $a_0 = 0$, $a_1 = 1$, $a_2 = -\frac{1}{2}$, $a_3 = \frac{1}{3}$, $a_4 = -\frac{1}{4}$, 当 $n=m=2$ 时,

由(7.11)得

$$-b_2 + \frac{1}{2}b_1 = \frac{1}{3},$$

$$\frac{1}{2}b_2 - \frac{1}{3}b_1 = -\frac{1}{4}.$$

求得 $b_1 = 1$, $b_2 = \frac{1}{6}$, 再由(7.9)得

$$a_0 = 0, a_1 = 1, a_2 = \frac{1}{2},$$

于是得

$$R_{2,2}(x) = \frac{x + \frac{1}{2}x^2}{1 + x + \frac{1}{6}x^2} = \frac{6x + 3x^2}{6 + 6x + x^2}.$$

当 $n=m=3$ 时, 由(7.11)得

$$-b_3 + \frac{1}{2}b_2 - \frac{1}{3}b_1 = -\frac{1}{4},$$

$$\frac{1}{2}b_3 - \frac{1}{3}b_2 + \frac{1}{4}b_1 = \frac{1}{5},$$

$$-\frac{1}{3}b_3 + \frac{1}{4}b_2 - \frac{1}{5}b_1 = -\frac{1}{6}.$$

解得 $b_1 = \frac{3}{2}$, $b_2 = \frac{3}{5}$, $b_3 = \frac{1}{20}$.

代入(7.9)得 $a_0 = 0$, $a_1 = 1$, $a_2 = 1$, $a_3 = \frac{11}{60}$.

于是得

$$R_{33}(x) = \frac{x + x^2 + \frac{11}{60}x^3}{1 + \frac{3}{2}x + \frac{3}{5}x^2 + \frac{1}{20}x^3} = \frac{60x + 60x^2 + 11x^3}{60 + 90x + 36x^2 + 3x^3}.$$

可以看到这里得到的 $R_{22}(x)$ 及 $R_{33}(x)$ 与 $\ln(1+x)$ 的前面连分式展开得到的有理逼近(7.4)结果一样.

为了求帕德逼近 $R_{nm}(x)$ 的误差估计, 由(7.9)及(7.11)求得的 $P_n(x), Q_m(x)$ 系数 a_0, a_1, \dots, a_n 及 b_0, b_1, \dots, b_m , 直接代入则得

$$f(x)Q_m(x) - P_n(x) = \sum_{l=0}^{m+n+1} b_k c_{n+m+1+l-k} x^l,$$

将 $Q_m(x)$ 除上式两端, 即得

$$f(x) - R_{nm}(x) = \frac{x^{n+m+1}}{\sum_{k=0}^{m+n+1} b_k c_{n+m+1+k} x^k}, \quad (7.13)$$

其中 $r = \sum_{k=0}^{m+n+1} b_k c_{n+m+1+k}$.

当 $|x| < 1$ 时可得误差近似表达式

$$f(x) - R_{nm}(x) = r x^{n+m+1}, \quad r = \sum_{k=0}^{m+n+1} b_k c_{n+m+1+k}.$$

评注

函数逼近是数学中的经典课题, 它与数学中其他分支有着密切的联系, 也是计算数学的基础, 本章仅讨论最佳一致逼近和最佳平方逼近的基本概念及正交多项式, 更详细内容可参看[5][6]等专门著作.

正交多项式在函数逼近中有重要作用, 它在高斯(Gauss)求积中也有重要应用. 特别地, 对勒让德多项式及切比雪夫多项式本章做了较详细的讨论, 因为这是两个十分重要又经常使用的正

交多项式, 应引起读者关注.

曲线拟合的最小二乘法在应用科学中具有重要作用, 它是离散点的最佳平方逼近. 本章引入哈尔条件, 由此可证明解的存在唯一性, 而采用离散点正交多项式可避免解法方程时出现的病态问题, 为用多项式作最小二乘模型提供了可行的算法.

傅里叶变换也是最佳平方逼近得到的, 快速傅里叶变换(FFT)是节省计算次数的一个范例, 本章介绍的 FFT 算法是 Cooley-Tukey 算法的改进^[12], 其基本思想与经典的 FFT 算法是一致的, 更详细讨论可参见文献[11].

有理逼近是函数逼近的重要组成部分. 本章只简单介绍有广泛应用的帕德逼近, 其他内容可参见文献[5].

习 题

1. $f(x) = \sin \frac{1}{2}x$, 给出 $[0, 1]$ 上的伯恩斯坦多项式 $B_1(f, x)$ 及 $B_3(f, x)$.

2. 当 $f(x) = x$ 时, 求证 $B_n(f, x) = x$.

3. 证明函数 $1, x, \dots, x^n$ 线性无关.

4. 计算下列函数 $f(x)$ 关于 $C[0, 1]$ 的 f , f^{-1} 与 f^{-2} :

$$(1) f(x) = (x - 1)^3,$$

$$(2) f(x) = \left| x - \frac{1}{2} \right|,$$

(3) $f(x) = x^m(1 - x)^n$, m 与 n 为正整数,

$$(4) f(x) = (x + 1)^{10} e^{-x}.$$

5. 证明 $f \cdot g = f^{-1} \circ g$.

6. 对 $f(x), g(x) \in C[a, b]$, 定义

$$(1) (f, g) = \int_a^b f(x)g(x)dx,$$

$$(2) (f, g) = \int_a^b f(x)g(x)dx + f(a)g(a).$$

问它们是否构成内积 .

7. 令 $T_n^*(x) = T_n(2x - 1)$, $x \in [0, 1]$, 试证 $\{T_n^*(x)\}$ 是在 $[0, 1]$ 上带权 $(x) = \frac{1}{x - x^2}$ 的正交多项式, 并求 $T_0^*(x), T_1^*(x), T_2^*(x), T_3^*(x)$.

8. 对权函数 $(x) = 1 + x^2$, 区间 $[-1, 1]$, 试求首项系数为 1 的正交多项式 $T_n(x)$, $n=0, 1, 2, 3$.

9. 试证明由(2.14)给出的第二类切比雪夫多项式族 $\{u_n(x)\}$ 是 $[-1, 1]$ 上带权 $(x) = 1 - x^2$ 的正交多项式 .

10. 证明切比雪夫多项式 $T_n(x)$ 满足微分方程

$$(1 - x^2) T_n'(x) - x T_n(x) + n^2 T_n(x) = 0.$$

11. 假设 $f(x)$ 在 $[a, b]$ 上连续, 求 $f(x)$ 的零次最佳一致逼近多项式 .

12. 选取常数 a , 使 $\max_{0 \leq x \leq 1} |x^3 - ax|$ 达到极小, 又问这个解是否唯一 ?

13. 求 $f(x) = \sin x$ 在 $[0, \pi/2]$ 上的最佳一次逼近多项式, 并估计误差 .

14. 求 $f(x) = e^x$ 在 $[0, 1]$ 上的最佳一次逼近多项式 .

15. 求 $f(x) = x^4 + 3x^3 - 1$ 在区间 $[0, 1]$ 上的三次最佳一致逼近多项式 .

16. $f(x) = |x|$, 在 $[-1, 1]$ 上求关于 $\mathcal{S} = \text{span}\{1, x^2, x^4\}$ 的最佳平方逼近多项式 .

17. 求函数 $f(x)$ 在指定区间上对于 $\mathcal{S} = \text{span}\{1, x\}$ 的最佳平方逼近多项式:

$$(1) f(x) = \frac{1}{x}, [1, 3]; \quad (2) f(x) = e^x, [0, 1];$$

$$(3) f(x) = \cos x, [0, 1]; \quad (4) f(x) = \ln x, [1, 2].$$

18. $f(x) = \sin \frac{\pi}{2} x$, 在 $[-1, 1]$ 上按勒让德多项式展开求三次最佳平方逼近多项式 .

19. 观测物体的直线运动, 得出以下数据:

时间 $t(s)$	0	0.9	1.9	3.0	3.9	5.0
距离 $s(m)$	0	10	30	50	80	110

求运动方程 .

20. 已知实验数据如下:

x_i	19	25	31	38	44
y_i	19. 0	32. 3	49. 0	73. 3	97. 8

用最小二乘法求形如 $y = a + bx^2$ 的经验公式，并计算均方误差。

21. 在某化学反应中，由实验得分解物浓度与时间关系如下：

时间 t	0	5	10	15	20	25	30	35	40	45	50	55
浓度 $y(\times 10^{-4})$	0	1. 27	2. 16	2. 86	3. 44	3. 87	4. 15	4. 37	4. 51	4. 58	4. 62	4. 64

用最小二乘法求 $y = f(t)$ 。

22. 给出一张记录 $\{f_k\} = (4, 3, 2, 1, 0, 1, 2, 3)$ ，用 FFT 算法求 $\{f_k\}$ 的离散谱 $\{\alpha\}$ 。

23. 用辗转相除法将 $R_{22}(x) = \frac{3x^2 + 6x}{x^2 + 6x + 6}$ 化为连分式。

24. 求 $f(x) = \sin x$ 在 $x = 0$ 处的(3, 3)阶帕德逼近 $R_{33}(x)$ 。

25. 求 $f(x) = e^x$ 在 $x = 0$ 处的(2, 1)阶帕德逼近 $R_{21}(x)$ 。

第4章 数值积分与数值微分

4.1 引言

4.1.1 数值求积的基本思想

实际问题当中常常需要计算积分。有些数值方法，如微分方程和积分方程的求解，也都和积分计算相联系。

依据人们所熟知的微积分基本定理，对于积分

$$I = \int_a^b f(x) dx,$$

只要找到被积函数 $f(x)$ 的原函数 $F(x)$ ，便有下列牛顿-莱布尼兹 (Newton-Leibniz) 公式：

$$\int_a^b f(x) dx = F(b) - F(a).$$

但实际使用这种求积方法往往有困难，因为大量的被积函数，诸如 $\frac{\sin x}{x}$, $\sin x^2$ 等等，找不到用初等函数表示的原函数；另外，当 $f(x)$ 是由测量或数值计算给出的一张数据表时，牛顿-莱布尼兹公式也不能直接运用。因此有必要研究积分的数值计算问题。

积分中值定理告诉我们，在积分区间 $[a, b]$ 内存在一点 c ，成立

$$\int_a^b f(x) dx = (b - a) f(c),$$

就是说，底为 $b - a$ 而高为 $f(c)$ 的矩形面积恰等于所求曲边梯形的面积 I (图 4-1)。问题在于点 c 的具体位置一般是不知道的，因而难以准确算出 $f(c)$ 的值。我们将 $f(c)$ 称为区间 $[a, b]$ 上的平均

高度。这样,只要对平均高度 $f(\cdot)$ 提供一种算法,相应地便获得一种数值求积方法。

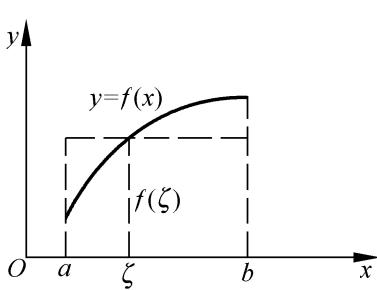


图 4-1

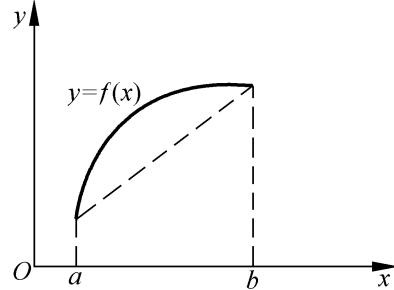


图 4-2

如果我们用两端点“高度” $f(a)$ 与 $f(b)$ 的算术平均作为平均高度 $f(\cdot)$ 的近似值,这样导出的求积公式

$$T = \frac{b-a}{2} [f(a) + f(b)] \quad (1.1)$$

便是我们所熟悉的梯形公式(几何意义参看图 4-2)。而如果改用区间中点 $c = \frac{a+b}{2}$ 的“高度” $f(c)$ 近似地取代平均高度 $f(\cdot)$,则又可导出所谓中矩形公式(今后简称矩形公式)

$$R = (b - a) f \left(\frac{a+b}{2} \right) . \quad (1.2)$$

更一般地,我们可以在区间 $[a, b]$ 上适当选取某些节点 x_k ,然后用 $f(x_k)$ 加权平均得到平均高度 $f(\cdot)$ 的近似值,这样构造出的求积公式具有下列形式:

$$\int_a^b f(x) dx = \sum_{k=0}^n A_k f(x_k), \quad (1.3)$$

式中 x_k 称为求积节点; A_k 称为求积系数,亦称伴随节点 x_k 的权。权 A_k 仅仅与节点 x_k 的选取有关,而不依赖于被积函数 $f(x)$ 的具体形式。

这类数值积分方法通常称为机械求积,其特点是将积分求值

问题归结为函数值的计算,这就避开了牛顿-莱布尼兹公式需要寻求原函数的困难.

4.1.2 代数精度的概念

数值求积方法是近似方法,为要保证精度,我们自然希望求积公式能对“尽可能多”的函数准确地成立,这就提出了所谓代数精度的概念.

定义 1 如果某个求积公式对于次数不超过 m 的多项式均能准确地成立,但对于 $m+1$ 次多项式就不准确成立,则称该求积公式具有 m 次代数精度.

不难验证,梯形公式(1.1)和矩形公式(1.2)均具有一次代数精度.

一般地,欲使求积公式(1.3)具有 m 次代数精度,只要令它对于 $f(x) = 1, x, \dots, x^m$ 都能准确成立,这就要求

$$\begin{aligned} A_k &= b - a, \\ A_k x_k &= \frac{1}{2}(b^2 - a^2), \\ &\dots\dots\dots \\ A_k x_k^m &= \frac{1}{m+1}(b^{m+1} - a^{m+1}). \end{aligned} \tag{1.4}$$

为简洁起见,这里省略了符号 $\underset{k=0}{\overset{n}{\dots}}$ 中的上下标.

如果我们事先选定求积节点 x_k ,譬如,以区间 $[a, b]$ 的等距分点作为节点,这时取 $m = n$ 求解方程组(1.4)即可确定求积系数 A_k ,而使求积公式(1.3)至少具有 n 次代数精度.本章第2节介绍这样一类求积公式,梯形公式是其中的一个特例.

为了构造出形如(1.3)的求积公式,原则上是一个确定参数 x_k 和 A_k 的代数问题.

4.1.3 插值型的求积公式

设给定一组节点

$$a = x_0 < x_1 < x_2 < \dots < x_n = b,$$

且已知函数 $f(x)$ 在这些节点上的值, 作插值函数 $L_n(x)$ (参看第 2 章(2.9)式). 由于代数多项式 $L_n(x)$ 的原函数是容易求出的, 我们取

$$I_n = \int_a^b L_n(x) dx$$

作为积分 $I = \int_a^b f(x) dx$ 的近似值, 这样构造出的求积公式

$$I_n = \sum_{k=0}^n A_k f(x_k) \quad (1.5)$$

称为是插值型的, 式中求积系数 A_k 通过插值基函数 $l_k(x)$ 积分得出

$$A_k = \int_a^b l_k(x) dx. \quad (1.6)$$

由插值余项定理(第 2 章的定理 2)即知, 对于插值型的求积公式(1.5), 其余项

$$R[f] = I - I_n = \int_a^b \frac{f^{(n+1)}(\xi)}{(n+1)!}(x) dx, \quad (1.7)$$

式中 与变量 x 有关, $\xi = (x - x_0)(x - x_1) \dots (x - x_n)$.

如果求积公式(1.5)是插值型的, 按式(1.7), 对于次数不超过 n 的多项式 $f(x)$, 其余项 $R[f]$ 等于零, 因而这时求积公式至少具有 n 次代数精度.

反之, 如果求积公式(1.5)至少具有 n 次代数精度, 则它必定是插值型的. 事实上, 这时公式(1.5)对于插值基函数 $l_k(x)$ 应准确成立, 即有

$$\int_a^b l_k(x) dx = \sum_{j=0}^n A_j l_k(x_j).$$

注意到 $l_k(x_j) = 1$, 上式右端实际上即等于 A_k , 因而式(1.6)

成立.

综上所述, 我们的结论是

定理1 形如(1.5)的求积公式至少有 n 次代数精度的充分必要条件是, 它是插值型的.

4.1.4 求积公式的收敛性与稳定性

定义2 在求积公式(1.3)中, 若

$$\lim_{\substack{n \\ h \rightarrow 0}} \sum_{k=0}^n A_k f(x_k) = \int_a^b f(x) dx.$$

其中 $h = \max_{1 \leq i \leq n} (x_i - x_{i-1})$, 则称求积公式(1.3)是收敛的.

在求积公式(1.3)中, 由于计算 $f(x_k)$ 可能产生误差 ε_k , 实际得到 \bar{f}_k , 即 $f(x_k) = \bar{f}_k + \varepsilon_k$. 记

$$I_n(f) = \sum_{k=0}^n A_k f(x_k), \quad I_n(\bar{f}) = \sum_{k=0}^n A_k \bar{f}_k.$$

如果对任给小正数 $\nu > 0$, 只要误差 $|\varepsilon_k|$ 充分小就有

$$|I_n(f) - I_n(\bar{f})| = \left| \sum_{k=0}^n A_k (f(x_k) - \bar{f}_k) \right|, \quad (1.8)$$

它表明求积公式(1.3)计算是稳定的, 由此给出:

定义3 对任给 $\nu > 0$, 若 $\forall k \in \{0, 1, \dots, n\}$ 有 $|f(x_k) - \bar{f}_k| < \nu$, 则称求积公式(1.3)是稳定的.

定理2 若求积公式(1.3)中系数 $A_k > 0$ ($k = 0, 1, \dots, n$), 则此求积公式是稳定的.

证明 对任给 $\nu > 0$, 若取 $\delta = \frac{\nu}{b-a}$, 对 $k = 0, 1, \dots, n$ 都有

$|f(x_k) - \bar{f}_k| < \delta$, 则有

$$\begin{aligned} |I_n(f) - I_n(\bar{f})| &= \left| \sum_{k=0}^n A_k (f(x_k) - \bar{f}_k) \right| \\ &\leq \sum_{k=0}^n |A_k| |f(x_k) - \bar{f}_k| \end{aligned}$$

n

$$\sum_{k=0}^n A_k = (b - a) = \dots$$

由定义 3 可知求积公式(1.3)是稳定的. 证毕.

定理 2 表明只要求积系数 $A_k > 0$, 就能保证计算的稳定性.

4.2 牛顿-柯特斯公式

4.2.1 柯特斯系数

设将积分区间 $[a, b]$ 划分为 n 等分, 步长 $h = \frac{b - a}{n}$, 选取等距

节点 $x_k = a + kh$ 构造出的插值型求积公式

$$I_n = (b - a) \sum_{k=0}^n C_k^{(n)} f(x_k) \quad (2.1)$$

称为牛顿-柯特斯 (Newton-Cotes) 公式, 式中 $C_k^{(n)}$ 称为柯特斯系数. 按(1.6)式, 引进变换 $x = a + th$, 则有

$$\begin{aligned} C_k^{(n)} &= \frac{h}{b - a} \int_0^n \frac{t - j}{k - j} dt \\ &= \frac{(-1)^{n-k}}{nk!(n-k)!} \int_0^n (t - j) dt. \end{aligned} \quad (2.2)$$

由于是多项式的积分, 柯特斯系数的计算不会遇到实质性的困难. 当 $n=1$ 时,

$$C_0^{(1)} = C_1^{(1)} = \frac{1}{2},$$

这时的求积公式就是我们所熟悉的梯形公式(1.1).

当 $n=2$ 时, 按(2.2)式, 这时柯特斯系数为

$$C_0^{(2)} = \frac{1}{4} \int_0^2 (t - 1)(t - 2) dt = \frac{1}{6},$$

$$C_1^{(2)} = -\frac{1}{2} \int_0^2 t(t - 2) dt = \frac{4}{6},$$

$$C_2^{(2)} = \frac{1}{4} \int_0^2 t(t-1) dt = \frac{1}{6}.$$

相应的求积公式是下列辛普森(Simpson)公式

$$S = \frac{b-a}{6} [f(a) + 4f\left(\frac{a+b}{2}\right) + f(b)], \quad (2.3)$$

而 $n=4$ 的牛顿-柯特斯公式则特别称为柯特斯公式, 其形式是

$$C = \frac{b-a}{90} [7f(x_0) + 32f(x_1) + 12f(x_2) + 32f(x_3) + 7f(x_4)]. \quad (2.4)$$

这里 $x_k = a + kh$, $h = \frac{b-a}{4}$.

下表列出柯特斯系数表开头的一部分.

表 4-1

n	$C_k^{(n)}$						
1	$\frac{1}{2}$	$\frac{1}{2}$					
2	$\frac{1}{6}$	$\frac{2}{3}$	$\frac{1}{6}$				
3	$\frac{1}{8}$	$\frac{3}{8}$	$\frac{3}{8}$	$\frac{1}{8}$			
4	$\frac{7}{90}$	$\frac{16}{45}$	$\frac{2}{15}$	$\frac{16}{45}$	$\frac{7}{90}$		
5	$\frac{19}{288}$	$\frac{25}{96}$	$\frac{25}{144}$	$\frac{25}{144}$	$\frac{25}{96}$	$\frac{19}{288}$	
6	$\frac{41}{840}$	$\frac{9}{35}$	$\frac{9}{280}$	$\frac{34}{105}$	$\frac{9}{280}$	$\frac{9}{35}$	$\frac{41}{840}$
7	$\frac{751}{17280}$	$\frac{3577}{17280}$	$\frac{1323}{17280}$	$\frac{2989}{17280}$	$\frac{2989}{17280}$	$\frac{1323}{17280}$	$\frac{3577}{17280}$
8	$\frac{989}{28350}$	$\frac{5888}{28350}$	$\frac{-928}{28350}$	$\frac{10496}{28350}$	$\frac{-4540}{28350}$	$\frac{10496}{28350}$	$\frac{-928}{28350}$

从表中看到 $n=8$ 时, 柯特斯系数 $C_k^{(n)}$ 出现负值, 于是有

$$\sum_{k=0}^n C_k^{(n)} > \sum_{k=0}^n C_k^{(n)} = 1,$$

特别地, 假定 $C_k^{(n)} (f(x_k) - \bar{f}_k) > 0$, 且 $|f(x_k) - \bar{f}_k| =$, 则有

$$\begin{aligned} |I_n(f) - I_n(\bar{f})| &= \left| \sum_{k=0}^n C_k^{(n)} [f(x_k) - \bar{f}_k] \right| \\ &= \sum_{k=0}^n C_k^{(n)} [f(x_k) - \bar{f}_k] \\ &= \sum_{k=0}^n |C_k^{(n)}| |f(x_k) - \bar{f}_k| = \sum_{k=0}^n |C_k^{(n)}| > . \end{aligned}$$

它表明初始数据误差将会引起计算结果误差增大, 即计算不稳定, 故 $n=8$ 的牛顿-柯特斯公式是不用的.

4.2.2 偶阶求积公式的代数精度

作为插值型的求积公式, n 阶的牛顿-柯特斯公式至少具有 n 次的代数精度(定理 1). 实际的代数精度能否进一步提高呢?

先看辛普森公式(2.3), 它是二阶牛顿-柯特斯公式, 因此至少具有二次代数精度. 进一步用 $f(x) = x^3$ 进行检验, 按辛普森公式计算得

$$S = \frac{b-a}{6} a^3 + 4 \frac{a+b}{2}^3 + b^3 .$$

另一方面, 直接求积得

$$I = \int_a^b x^3 dx = \frac{b^4 - a^4}{4} .$$

这时有 $S = I$, 即辛普森公式对次数不超过三次的多项式均能准确成立, 又容易验证它对 $f(x) = x^4$ 通常是不准确的, 因此, 辛普森公式实际上具有三次代数精度.

一般地, 我们可以证明下述论断:

定理 3 当阶 n 为偶数时, 牛顿-柯特斯公式(2.1)至少有 $n+1$ 次代数精度.

证明 我们只要验证, 当 n 为偶数时, 牛顿-柯特斯公式对 $f(x) = x^{n+1}$ 的余项为零.

按余项公式(1.7), 由于这里 $f^{(n+1)}(x) = (n+1)!$, 从而有

$$R[f] = \int_a^b (x - x_j)^n dx,$$

引进变换 $x = a + th$, 并注意到 $x_j = a + jh$, 有

$$R[f] = h^{n+2} \int_0^n (t - j)^n dt,$$

若 n 为偶数, 则 $\frac{n}{2}$ 为整数, 再令 $t = u + \frac{n}{2}$, 进一步有

$$R[f] = h^{n+2} \int_{-\frac{n}{2}}^{\frac{n}{2}} u + \frac{n}{2} - j du,$$

据此可以断定 $R[f] = 0$, 因为被积函数

$$H(u) = \int_{j=0}^n u + \frac{n}{2} - j = \int_{j=-\frac{n}{2}}^{\frac{n}{2}} (u - j)$$

是个奇函数. 证毕.

4.2.3 几种低阶求积公式的余项

首先考察梯形公式, 按余项公式(1.7), 梯形公式(1.1)的余项

$$R_T = I - T = \int_a^b \frac{f(\xi)}{2}(x - a)(x - b) dx,$$

这里积分的核函数 $(x - a)(x - b)$ 在区间 $[a, b]$ 上保号(非正), 应用积分中值定理, 在 $[a, b]$ 内存在一点 ξ , 使

$$\begin{aligned} R_T &= \frac{f(\xi)}{2} \int_a^b (x - a)(x - b) dx \\ &= -\frac{f(\xi)}{12}(b - a)^3, \quad [\alpha, \beta]. \end{aligned} \tag{2.5}$$

再研究辛普森公式(2.3)的余项 $R_S = I - S$. 为此构造次数不超过 3 的多项式 $H(x)$, 使满足

$$H(a) = f(a), \quad H(b) = f(b);$$

$$H(c) = f(c), \quad H'(c) = f'(c). \quad (2.6)$$

这里 $c = \frac{a+b}{2}$. 由于辛普森公式具有三次代数精度, 它对于这样构造出的三次式 $H(x)$ 是准确的, 即

$$\int_a^b H(x) dx = \frac{b-a}{6} [H(a) + 4H(c) + H(b)],$$

而利用插值条件 (2.6) 知, 上式右端实际上等于按辛普森公式 (2.3) 求得的积分值 S , 因此积分余项

$$R_S = I - S = \int_a^b [f(x) - H(x)] dx.$$

对于满足条件 (2.6) 的多项式 $H(x)$, 其插值余项由第 2 章 (5.11) 得

$$f(x) - H(x) = \frac{f^{(4)}(\cdot)}{4!}(x-a)(x-c)^2(x-b),$$

故有

$$R_S = \int_a^b \frac{f^{(4)}(\cdot)}{4!}(x-a)(x-c)^2(x-b) dx.$$

这时积分的核函数 $(x-a)(x-c)^2(x-b)$ 在 $[a, b]$ 上保号(非正), 再用积分中值定理有

$$\begin{aligned} R_S &= \frac{f^{(4)}(\cdot)}{4!} \int_a^b (x-a)(x-c)^2(x-b) dx \\ &= -\frac{b-a}{180} \frac{b-a}{2} f^{(4)}(\cdot). \end{aligned} \quad (2.7)$$

关于柯特斯公式 (2.4) 的积分余项, 这里不再具体推导, 仅列出结果如下:

$$R_C = I - C = -\frac{2(b-a)}{945} \frac{b-a}{4} f^{(6)}(\cdot). \quad (2.8)$$

4.3 复化求积公式

前面已经指出高阶牛顿-柯特斯公式是不稳定的, 因此, 不可

能通过提高阶的方法来提高求积精度.为了提高精度通常可把积分区间分成若干子区间(通常是等分),再在每个子区间上用低阶求积公式.这种方法称为复化求积法.本节只讨论复化梯形公式与复化辛普森公式.

4.3.1 复化梯形公式

将区间 $[a, b]$ 划分为 n 等分,分点 $x_k = kh$, $h = \frac{b-a}{n}$, $k=0, 1, \dots, n$,在每个子区间 $[x_k, x_{k+1}]$ ($k=0, 1, \dots, n-1$)上采用梯形公式(1.1),则得

$$\begin{aligned} I &= \int_a^b f(x) dx = \sum_{k=0}^{n-1} \int_{x_k}^{x_{k+1}} f(x) dx \\ &= \frac{h}{2} \sum_{k=0}^{n-1} [f(x_k) + f(x_{k+1})] + R_n(f). \quad (3.1) \end{aligned}$$

记

$$\begin{aligned} T_n &= \frac{h}{2} \sum_{k=0}^{n-1} [f(x_k) + f(x_{k+1})] \\ &= \frac{h}{2} f(a) + 2 \sum_{k=1}^{n-1} f(x_k) + f(b), \quad (3.2) \end{aligned}$$

称为复化梯形公式,其余项可由(2.5)得

$$R_n(f) = I - T_n = \sum_{k=0}^{n-1} -\frac{h^3}{12} f''(x_k), \quad x_k \in (x_k, x_{k+1}).$$

由于 $f(x) \in C^2[a, b]$,且

$$0 \min_{k=n-1} f'(x_k) - \frac{1}{n} \sum_{k=0}^{n-1} f'(x_k) \max_{k=n-1} f'(x_k).$$

所以 $v(a, b)$ 使

$$f(\) = \frac{1}{n} \sum_{k=0}^{n-1} f'(x_k).$$

于是复化梯形公式余项为

$$R_n(f) = -\frac{b-a}{12}h^2 f''(\xi). \quad (3.3)$$

可以看出误差是 h^2 阶, 且由 (3.3) 立即得到, 当 $f(x) \in C[a, b]$, 则

$$\lim_{n \rightarrow \infty} T_n = \int_a^b f(x) dx,$$

即复化梯形公式是收敛的. 事实上只要设 $f(x) \in C[a, b]$, 则可得到收敛性, 因为只要把 T_n 改写为

$$T_n = \frac{1}{2} \left[\frac{b-a}{n} \sum_{k=0}^{n-1} f(x_k) + \frac{b-a}{n} \sum_{k=1}^n f(x_k) \right].$$

当 n 时, 上式右端括号内的两个和式均收敛到积分 $\int_a^b f(x) dx$, 所以复化梯形公式 (3.2) 收敛. 此外, T_n 的求积系数为正, 由定理 2 知复化梯形公式是稳定的.

4.3.2 复化辛普森求积公式

将区间 $[a, b]$ 分为 n 等分, 在每个子区间 $[x_k, x_{k+1}]$ 上采用辛普森公式 (2.3), 若记 $x_{k+1/2} = x_k + \frac{1}{2}h$, 则得

$$\begin{aligned} I &= \int_a^b f(x) dx = \sum_{k=0}^{n-1} \int_{x_k}^{x_{k+1}} f(x) dx \\ &= \frac{h}{6} \sum_{k=0}^{n-1} [f(x_k) + 4f(x_{k+1/2}) + f(x_{k+1})] + R_n(f). \end{aligned} \quad (3.4)$$

记

$$\begin{aligned} S_n &= \frac{h}{6} \sum_{k=0}^{n-1} [f(x_k) + 4f(x_{k+1/2}) + f(x_{k+1})] \\ &= \frac{h}{6} [f(a) + 4 \sum_{k=0}^{n-1} f(x_{k+1/2}) + 2 \sum_{k=1}^{n-1} f(x_k) + f(b)], \end{aligned} \quad (3.5)$$

称为复化辛普森求积公式. 其余项由 (2.7) 得

$$R_n(f) = I - S_n = -\frac{h}{180} \sum_{k=0}^{n-1} \frac{h}{2} f^{(4)}(x_k), \quad k \in (x_k, x_{k+1}),$$

于是当 $f(x) \in C^4[a, b]$ 时, 与复化梯形公式相似有

$$R_n(f) = I - S_n = -\frac{b-a}{180} \sum_{k=0}^{n-1} \frac{h}{2} f^{(4)}(x_k), \quad (a, b). \quad (3.6)$$

由(3.6)看出, 误差阶为 h^4 , 收敛性是显然的, 实际上, 只要 $f(x) \in C[a, b]$ 则可得到收敛性, 即

$$\lim_{n \rightarrow \infty} S_n = \int_a^b f(x) dx.$$

此外, 由于 S_n 中求积系数均为正数, 故知复化辛普森公式计算稳定.

例 1 对于函数 $f(x) = \frac{\sin x}{x}$, 给出 $n = 8$ 的函数表(见表 4-2), 试用复化梯形公式(3.2)及复化辛普森公式(3.5)计算积分

$$I = \int_0^1 \frac{\sin x}{x} dx,$$

并估计误差.

解 将积分区间 $[0, 1]$ 划分为 8 等分, 应用复化梯形法求得

$$T_8 = 0.9456909;$$

而如果将 $[0, 1]$ 分为 4 等分, 应用复化辛普森法有

$$S_4 = 0.9460832.$$

比较上面两个结果 T_8 与 S_4 , 它们都需要提供 9 个点上的函数值, 计算量基本相同, 然而精度却差别很大, 同积分的准确值 $I = 0.9460831$ 比较, 复化梯形法的结果

表 4-2

x	f(x)
0	1
1/8	0.9973978
1/4	0.9896158
3/8	0.9767267
1/2	0.9588510
5/8	0.9361556
3/4	0.9088516
7/8	0.8771925
1	0.8414709

果 $T_8 = 0.9456909$ 只有两位有效数字, 而复化辛普森法的结果 $S_4 = 0.9460832$ 却有六位有效数字 .

为了利用余项公式估计误差, 要求 $f(x) = \frac{\sin x}{x}$ 的高阶导数,

由于

$$f(x) = \frac{\sin x}{x} = \int_0^1 \cos(xt) dt,$$

所以有

$$f^{(k)}(x) = \int_0^1 \frac{d^k}{dx^k} (\cos xt) dt = \int_0^1 t^k \cos xt + \frac{k}{2} dt,$$

于是

$$\max_{0 \leq x \leq 1} |f^{(k)}(x)| = \left| \int_0^1 \cos xt + \frac{k}{2} dt \right| = \int_0^1 t^k dt = \frac{1}{k+1}.$$

由(3.3)得复化梯形公式误差

$$|R_8(f)| = |I - T_n| = \frac{h^2}{12} \max_{0 \leq x \leq 1} |f'(x)| \\ = \frac{1}{12} \cdot \frac{1}{8}^2 \cdot \frac{1}{3} = 0.000434.$$

对复化辛普森公式误差, 由(3.6)得

$$|R_4(f)| = |I - S_4| = \frac{1}{2880} \cdot \frac{1}{4}^4 \cdot \frac{1}{5} = 0.271 \times 10^{-6}.$$

4.4 龙贝格求积公式

4.4.1 梯形法的递推化

上节介绍的复化求积方法可提高求积精度, 实际计算时若精度不够可将步长逐次分半. 设将区间 $[a, b]$ 分为 n 等分, 共有 $n+1$ 个分点, 如果将求积区间再二分一次, 则分点增至 $2n+1$ 个, 我们将二分前后两个积分值联系起来加以考察. 注意到每个子区间

$[x_k, x_{k+1}]$ 经过二分只增加了一个分点 $x_{k+\frac{1}{2}} = \frac{1}{2}(x_k + x_{k+1})$, 用复化梯形公式求得该子区间上的积分值为

$$\frac{h}{4} [f(x_k) + 2f(x_{k+\frac{1}{2}}) + f(x_{k+1})].$$

注意, 这里 $h = \frac{b-a}{n}$ 代表二分前的步长. 将每个子区间上的积分值相加得

$$T_{2n} = \frac{h}{4} \sum_{k=0}^{n-1} [f(x_k) + f(x_{k+1})] + \frac{h}{2} \sum_{k=0}^{n-1} f(x_{k+\frac{1}{2}}),$$

从而利用式(3.2)可导出下列递推公式

$$T_{2n} = \frac{1}{2} T_n + \frac{h}{2} \sum_{k=0}^{n-1} f(x_{k+\frac{1}{2}}). \quad (4.1)$$

例2 计算积分值

$$I = \int_0^1 \frac{\sin x}{x} dx.$$

解 我们先对整个区间 $[0, 1]$ 使用梯形公式. 对于函数 $f(x) = \frac{\sin x}{x}$, 它在 $x=0$ 的值定义为 $f(0) = 1$, 而 $f(1) = 0.8414709$, 据梯形公式计算得

$$T_1 = \frac{1}{2} [f(0) + f(1)] = 0.9207355.$$

然后将区间二等分, 再求出中点的函数值

$$f\left(\frac{1}{2}\right) = 0.9588510,$$

从而利用递推公式(4.1), 有

$$T_2 = \frac{1}{2} T_1 + \frac{1}{2} f\left(\frac{1}{2}\right) = 0.9397933.$$

我们进一步二分求积区间, 并计算新分点上的函数值

$$f(1/4) = 0.9896158, \quad f(3/4) = 0.9088516.$$

再利用式(4.1), 有

$$T_4 = \frac{1}{2} T_2 + \frac{1}{4} f\left(\frac{1}{4}\right) + f\left(\frac{3}{4}\right) = 0.9445135.$$

这样不断二分下去, 计算结果见下表(表中 k 代表二分次数, 区间等分数 $n = 2^k$).

表 4-3

k	1	2	3	4	5
T_n	0.9397933	0.9445135	0.9456909	0.9459850	0.9460596
k	6	7	8	9	10
T_n	0.9460769	0.9460815	0.9460827	0.9460830	0.9460831

它表明用复化梯形公式计算积分 I 要达到 7 位有效数字的精度需要二分区间 10 次, 即要有分点 1025 个, 计算量很大.

4.4.2 龙贝格算法

梯形法计算简单但收敛慢, 如何提高收敛速度以节省计算量是本节要讨论的中心问题. 根据复化梯形公式的余项表达式(3.3)可知

$$I - T_n = -\frac{b-a}{12} h^2 f'(\xi), \quad (a, b);$$

$$I - T_{2n} = -\frac{b-a}{12} \frac{h}{2}^2 f''(\xi), \quad (\xi \in (a, b)).$$

假定 $f'(\xi) = f''(\xi)$, 则有

$$\frac{I - T_{2n}}{I - T_n} = \frac{1}{4}.$$

将上式移项整理, 可得

$$I - T_{2n} = \frac{1}{3} (T_{2n} - T_n). \quad (4.2)$$

由此可见, 只有二分前后的两个积分值 T_n 与 T_{2n} 相当接近, 就可以保证计算结果 T_{2n} 的误差很小. 这样直接用计算结果来估计误差的方法通常称作误差的事后估计法.

按式(4.2), 积分近似值 T_{2n} 的误差大致等于 $\frac{1}{3}(T_{2n} - T_n)$, 因此如果用这个误差值作为 T_{2n} 的一种补偿, 可以期望, 所得到的

$$\bar{T} = T_{2n} + \frac{1}{3}(T_{2n} - T_n) = \frac{4}{3}T_{2n} - \frac{1}{3}T_n \quad (4.3)$$

可能是更好的结果.

再考察例 2, 所求得的两个梯形值 $T_4 = 0.9445135$ 和 $T_8 = 0.9456909$ 的精度都很差(与准确值 $I = 0.9460831$ 比较, 只有两、三位有效数字), 但如果将它们按式(4.3)做线性组合, 则新的近似值

$$\bar{T} = \frac{4}{3}T_8 - \frac{1}{3}T_4 = 0.9460833$$

却有 6 位有效数字.

按公式(4.3)组合得到的近似值 \bar{T} 其实质究竟是什么呢? 直接验证易知

$$S_n = \frac{4}{3}T_{2n} - \frac{1}{3}T_n. \quad (4.4)$$

这就是说, 用梯形法二分前后的两个积分值 T_n 与 T_{2n} , 按式(4.3)做线性组合, 结果得到辛普森法的积分值 S_n .

再考察辛普森法, 按误差公式(3.6), 其截断误差大致与 h^4 成正比, 因此, 若将步长折半则误差将减至原有误差的 $1/16$, 即有

$$\frac{I - S_{2n}}{I - S_n} = \frac{1}{16},$$

由此得到

$$I = \frac{16}{15}S_{2n} - \frac{1}{15}S_n.$$

不难直接验证, 上式右端的值等于 C_n , C_n 为复化柯特斯公式, 它的精度为 $O(h^6)$, 就是说, 用辛普森法二分前后两个积分值 S_n 与 S_{2n} , 可得到 S_n 的近似误差估计为 $\frac{1}{15}(S_{2n} - S_n)$, 并且按上式做线性组合可得复化柯特斯的积分值 C_n , 即有

$$C_n = \frac{16}{15}S_{2n} - \frac{1}{15}S_n. \quad (4.5)$$

重复同样的手续, 依据柯特斯法的误差阶为 h^6 , 可进一步导出下列龙贝格(Romberg)公式:

$$R_n = \frac{64}{63}C_{2n} - \frac{1}{63}C_n. \quad (4.6)$$

我们在变步长的过程中运用公式(4.4), (4.5)和(4.6), 就能将粗糙的梯形值 T_n 逐步加工成精度较高的辛普森值 S_n 、柯特斯值 C_n 和龙贝格值 R_n .

例 3 用加速公式(4.4), (4.5)和(4.6)加工例 2 得到的梯形值, 计算结果见下表(k 代表二分次数):

表 4-4

k	T_{2^k}	$S_{2^{k-1}}$	$C_{2^{k-2}}$	$R_{2^{k-3}}$
0	0.9207355			
1	0.9397933	0.9461459		
2	0.9445135	0.9460869	0.9460830	
3	0.9456909	0.9460833	0.9460831	0.9460831

我们看到, 这里利用二分 3 次的数据(它们的精度都很差, 只有两三位有效数字), 通过三次加速求得 $R = 0.9460831$, 这个结果的每一位数字都是有效数字, 可见加速的效果是十分显著的.

4.4.3 理查森外推加速法

上面讨论说明由梯形公式出发, 将区间 $[a, b]$ 逐次二分可提高

求积公式的精度, 上述加速过程还可继续下去, 其理论依据是梯形公式的余项展开, 设

$$I - T_n = -\frac{b-a}{12}h^2 f''(\xi), \quad [a, b], \quad h = \frac{b-a}{n}.$$

若记 $T_n = T(h)$, 当区间 $[a, b]$ 分为 $2n$ 等分时, 则有 $T_{2n} = T\left(\frac{h}{2}\right)$,

并且有

$$T(h) = I + \frac{b-a}{12}h^2 f''(\xi), \quad \lim_{h \rightarrow 0} T(h) = T(0) = I,$$

可证明梯形公式余项可展成级数形式, 即

定理4 设 $f(x) \in C[a, b]$, 则有

$$T(h) = I + r_1 h^2 + r_2 h^4 + \dots + r_l h^{2l} + \dots, \quad (4.7)$$

其中系数 r_l ($l=1, 2, \dots$) 与 h 无关.

此定理可利用 $f(x)$ 的泰勒展开推导得到, 此处从略.

定理4 表明 $T(h) - I$ 是 $O(h^2)$ 阶, 在 (4.7) 中, 若用 $h/2$ 代替 h , 有

$$T\left(\frac{h}{2}\right) = I + r_1 \frac{h^2}{4} + r_2 \frac{h^4}{16} + \dots + r_l \frac{h^{2l}}{2^{2l}} + \dots, \quad (4.8)$$

若用 4 乘 (4.8) 式, 减去 (4.7) 式再除 3 记之为 $T_1(h)$, 则得

$$T_1(h) = \frac{4T\left(\frac{h}{2}\right) - T(h)}{3} = I + r_1 h^4 + r_2 h^6 + \dots, \quad (4.9)$$

这里 r_1, r_2, \dots 以及后面将出现的 r_k, r_k 均为与 h 无关的系数, 这样构造的 $T_1(h)$ 与积分值 I 近似的阶为 $O(h^4)$. 比较 (4.9) 与 (4.4)

可知, 这样构造的序列 $T_1(h), T_1\left(\frac{h}{2}\right), \dots$ 就是辛普森公式序列

S_n, S_{2n}, \dots .

又根据 (4.9) 有

$$T_1 - \frac{h}{2} = I + {}_1 \frac{h^4}{16} + {}_2 h^6 + {}_3 h^8 + \dots,$$

若令

$$T_2(h) = \frac{16}{15} T_1 - \frac{h}{2} - \frac{1}{16} T_1(h),$$

则又可进一步从余项展开式中消去 h^4 项, 而有

$$T_2(h) = I + {}_1 h^6 + {}_2 h^8 + \dots.$$

这样构造出的 $\{T_2(h)\}$, 其实就是柯特斯公式序列, 它与积分值 I 的逼近阶为 $O(h^6)$. 如此继续下去, 每加速一次, 误差的量级便提高 2 阶, 一般地, 若记 $T_0(h) = T(h)$, 则有

$$T_m(h) = \frac{4^m}{4^m - 1} T_{m-1} - \frac{h}{2} - \frac{1}{4^m - 1} T_{m-1}(h). \quad (4.10)$$

经过 m ($m=1, 2, \dots$) 次加速后, 余项便取下列形式:

$$T_m(h) = I + {}_1 h^{2(m+1)} + {}_2 h^{2(m+2)} + \dots. \quad (4.11)$$

上述处理方法通常称为理查森(Richardson)外推加速方法.

设以 $T_0^{(k)}$ 表示二分 k 次后求得的梯形值, 且以 $T_m^{(k)}$ 表示序列 $\{T_0^{(k)}\}$ 的 m 次加速值, 则依递推公式(4.10)可得

$$T_m^{(k)} = \frac{4^m}{4^m - 1} T_{m-1}^{(k+1)} - \frac{1}{4^m - 1} T_{m-1}^{(k)} \quad (k = 1, 2, \dots). \quad (4.12)$$

公式(4.12)也称为龙贝格求积算法, 计算过程如下:

(1) 取 $k=0$, $h=b-a$, 求 $T_0^{(0)} = \frac{h}{2}[f(a) + f(b)]$.

令 1 k (记区间 $[a, b]$ 的二分次数).

(2) 求梯形值 $T_0 = \frac{b-a}{2^k}$, 即按递推公式(4.1)计算 $T_0^{(k)}$.

(3) 求加速值, 按公式(4.12)逐个求出 T 表(见表 4-5)的第 k 行其余各元素 $T_j^{(k-j)}$ ($j=1, 2, \dots, k$).

(4) 若 $|T_k^{(0)} - T_{k-1}^{(0)}| <$ (预先给定的精度), 则终止计算, 并取

$T_k^{(0)}$ I ; 否则令 $k+1 = k$ 转(2)继续计算 .

表 4-5 \mathbf{T} 表

k	h	$T_0^{(k)}$	$T_1^{(k)}$	$T_2^{(k)}$	$T_3^{(k)}$	$T_4^{(k)}$...
0	$b - a$	$T_0^{(0)}$					
1	$\frac{b-a}{2}$	$T_0^{(1)}$	$T_1^{(0)}$				
2	$\frac{b-a}{4}$	$T_0^{(2)}$	$T_1^{(1)}$	$T_2^{(0)}$			
3	$\frac{b-a}{8}$	$T_0^{(3)}$	$T_1^{(2)}$	$T_2^{(1)}$	$T_3^{(0)}$		
4	$\frac{b-a}{16}$	$T_0^{(4)}$	$T_1^{(3)}$	$T_2^{(2)}$	$T_3^{(1)}$	$T_4^{(0)}$	
...	W

表 4-5 指出了计算过程, 第 2 列 $h = \frac{b-a}{2^k}$ 给出了子区间长度, i 表示第 i 步外推 .

可以证明, 如果 $f(x)$ 充分光滑, 那么 T 表每一列的元素及对角线元素均收敛到所求的积分值 I , 即

$$\lim_k T_m^{(k)} = I \quad (m \text{ 固定}),$$

$$\lim_m T_m^{(0)} = I.$$

对于 $f(x)$ 不充分光滑的函数也可用龙贝格算法计算, 只是收敛慢一些, 这时也可以直接使用复化辛普森公式计算 . 见下面例题 .

例 4 用龙贝格算法计算积分 $I = \int_0^1 x^{3/2} dx$.

解 $f(x) = x^{3/2}$ 在 $[0, 1]$ 上仅是一次连续可微, 用龙贝格算法计算结果见表 4-6 . 从表中看到用龙贝格算到 $k=5$ 的精度与辛普森求积精度相当 . 这里 I 的精确值为 0.4 .

表 4.6

k	$T_0^{(k)}$	$T_1^{(k)}$	$T_2^{(k)}$	$T_3^{(k)}$	$T_4^{(k)}$	$T_5^{(k)}$
0	0.500000					
1	0.426777	0.402369				
2	0.407018	0.400432	0.400302			
3	0.401812	0.400077	0.400054	0.400050		
4	0.400463	0.400014	0.400009	0.400009	0.400009	
5	0.400118	0.400002	0.400002	0.400002	0.400002	0.400002

4.5 高斯求积公式

4.5.1 一般理论

形如(1.3)的机械求积公式

$$\int_a^b f(x) dx = \sum_{k=0}^n A_k f(x_k)$$

含有 $2n+2$ 个待定参数 $x_k, A_k (k=0, 1, \dots, n)$. 当 x_k 为等距节点时得到的插值求积公式其代数精度至少为 n 次, 如果适当选取 $x_k (k=0, 1, \dots, n)$, 有可能使求积公式具有 $2n+1$ 次代数精度, 这类求积公式称为高斯(Gauss)求积公式. 为使问题更具一般性, 我们

研究带权积分 $I = \int_a^b f(x) \omega(x) dx$, 这里 $\omega(x)$ 为权函数, 类似

(1.3), 它的求积公式为

$$\int_a^b f(x) \omega(x) dx = \sum_{k=0}^n A_k f(x_k), \quad (5.1)$$

$A_k (k=0, 1, \dots, n)$ 为不依赖于 $f(x)$ 的求积系数, $x_k (k=0, 1, \dots, n)$ 为求积节点, 可适当选取 x_k 及 $A_k (k=0, 1, \dots, n)$ 使(5.1)具有 $2n+1$ 次代数精度.

定义4 如果求积公式(5.1)具有 $2n+1$ 次代数精度, 则称其节点 $x_k (k=0, 1, \dots, n)$ 为高斯点, 相应公式(5.1)称为高斯求积公式.

根据定义要使(5.1)具有 $2n+1$ 次代数精度, 只要取 $f(x)=x^m$, 对 $m=0, 1, \dots, 2n+1$, (5.1)精确成立, 则得

$$\sum_{k=0}^n A_k x_k = \int_a^b x^m (x) dx \quad m = 0, 1, \dots, 2n+1. \quad (5.2)$$

当给定权函数 (x) , 求出右端积分, 则可由(5.2)解得 A_k 及 $x_k (k=0, 1, \dots, n)$.

例5 试构造下列积分的高斯求积公式:

$$\int_0^1 x f(x) dx = A_0 f(x_0) + A_1 f(x_1). \quad (5.3)$$

解 令公式(5.3)对于 $f(x)=1, x, x^2, x^3$ 准确成立, 得

$$\begin{aligned} A_0 + A_1 &= \frac{2}{3}; \\ x_0 A_0 + x_1 A_1 &= \frac{2}{5}; \\ x_0^2 A_0 + x_1^2 A_1 &= \frac{2}{7}; \\ x_0^3 A_0 + x_1^3 A_1 &= \frac{2}{9}. \end{aligned} \quad (5.4)$$

由于

$$x_0 A_0 + x_1 A_1 = x_0 (A_0 + A_1) + (x_1 - x_0) A_1,$$

利用(5.4)的第1式, 可将第2式化为

$$\frac{2}{3} x_0 + (x_1 - x_0) A_1 = \frac{2}{5}.$$

同样地, 利用第2式化第3式, 利用第3式化第4式, 分别得

$$\frac{2}{5} x_0 + (x_1 - x_0) x_1 A_1 = \frac{2}{7};$$

$$\frac{2}{7}x_0 + (x_1 - x_0)x_1^2 A_1 = \frac{2}{9}.$$

从上面三个式子消去 $(x_1 - x_0)A_1$, 有

$$\frac{2}{5}x_0 + \frac{2}{5} - \frac{2}{3}x_0 x_1 = \frac{2}{7};$$

$$\frac{2}{7}x_0 + \frac{2}{7} - \frac{2}{5}x_0 x_1 = \frac{2}{9}.$$

进一步整理得

$$\frac{2}{5}(x_0 + x_1) - \frac{2}{3}x_0 x_1 = \frac{2}{7};$$

$$\frac{2}{7}(x_0 + x_1) - \frac{2}{5}x_0 x_1 = \frac{2}{9}.$$

由此解出

$$x_0 x_1 = \frac{5}{21}, \quad x_0 + x_1 = \frac{10}{9},$$

从而求出

$$x_0 = 0.821162, \quad x_1 = 0.289949;$$

$$A_0 = 0.389111, \quad A_1 = 0.277556.$$

于是形如(5.3)的高斯公式是

$$\begin{aligned} & \int_0^1 x f(x) dx = 0.389111 f(0.821162) \\ & \quad + 0.277556 f(0.289949). \end{aligned}$$

从此例看到求解非线性方程组(5.2)较复杂, 通常 $n \geq 2$ 就很难求解. 故一般不通过解方程(5.2)求 x_k 及 A_k ($k = 0, 1, \dots, n$), 而从分析高斯点的特性来构造高斯求积公式.

定理5 插值型求积公式(5.1)的节点 $a = x_0 < x_1 < \dots < x_n$ 是高斯点的充分必要条件是以这些节点为零点的多项式

$$n+1(x) = (x - x_0)(x - x_1) \dots (x - x_n)$$

与任何次数不超过 n 的多项式 $P(x)$ 带权 $\varphi(x)$ 正交, 即

$$\int_a^b P(x)_{n+1}(x) f(x) dx = 0. \quad (5.5)$$

证明 必要性. 设 $P(x) \in H_n$, 则 $P(x)_{n+1}(x) \in H_{2n+1}$, 因此, 如果 x_0, x_1, \dots, x_n 是高斯点, 则求积公式(5.1)对于 $f(x) = P(x)_{n+1}(x)$ 精确成立, 即有

$$\int_a^b P(x)_{n+1}(x) f(x) dx = \sum_{k=0}^n A_k P(x_k)_{n+1}(x_k).$$

因 $P(x)_{n+1}(x_k) = 0 (k=0, 1, \dots, n)$, 故(5.5)成立.

再证充分性. 对于 " $f(x) \in H_{2n+1}$ ", 用 $P(x)_{n+1}(x)$ 除 $f(x)$, 记商为 $P(x)$, 余式为 $q(x)$, 即 $f(x) = P(x)_{n+1}(x) + q(x)$, 其中 $P(x), q(x) \in H_n$. 由(5.5)可得

$$\int_a^b f(x) f(x) dx = \int_a^b q(x) f(x) dx. \quad (5.6)$$

由于所给求积公式(5.1)是插值型的, 它对于 $q(x) \in H_n$ 是精确的, 即

$$\int_a^b q(x) f(x) dx = \sum_{k=0}^n A_k q(x_k).$$

再注意到 $q(x_k) = 0 (k=0, 1, \dots, n)$, 知 $q(x_k) = f(x_k) (k=0, 1, \dots, n)$, 从而由(5.6)有

$$\int_a^b f(x) f(x) dx = \int_a^b q(x) f(x) dx = \sum_{k=0}^n A_k f(x_k).$$

可见求积公式(5.1)对一切次数不超过 $2n+1$ 的多项式均精确成立. 因此, $x_k (k=0, 1, \dots, n)$ 为高斯点. 证毕.

定理表明在 $[a, b]$ 上带权 (x) 的 $n+1$ 次正交多项式的零点就是求积公式(5.1)的高斯点, 有了求积节点 $x_k (k=0, 1, \dots, n)$, 再利用(5.2)对 $m=0, 1, \dots, n$ 成立, 则得到一组关于求积系数 A_0, A_1, \dots, A_n 的线性方程. 解此方程则得 $A_k (k=0, 1, \dots, n)$. 也可直接由 x_0, x_1, \dots, x_n 的插值多项式求出求积系数 $A_k (k=0, 1, \dots, n)$.

$0, 1, \dots, n)$.

下面讨论高斯求积公式(5.1)的余项. 利用 $f(x)$ 在节点 x_k ($k = 0, 1, \dots, n$) 的埃尔米特插值 $H_{2n+1}(x)$, 即

$$H_{2n+1}(x_k) = f(x_k), \quad H_{2n+1}(x_k) = f(x_k), \quad k = 0, 1, \dots, n.$$

于是

$$f(x) = H_{2n+1}(x) + \frac{f^{(2n+2)}(\cdot)}{(2n+2)!} {}_{n+1}^2(x)$$

两端乘 (x) , 并由 a 到 b 积分, 则得

$$I = \int_a^b f(x)(x) dx = \int_a^b H_{2n+1}(x)(x) dx + R_n[f]. \quad (5.7)$$

其中右端第一项积分对 $2n+1$ 次多项式精确成立, 故

$$R_n[f] = I - \sum_{k=0}^n A_k f(x_k) = \int_a^b \frac{f^{(2n+2)}(\cdot)}{(2n+2)!} {}_{n+1}^2(x)(x) dx.$$

由于 ${}_{n+1}^2(x)(x) \equiv 0$, 故由积分中值定理得(5.1)的余项为

$$R_n[f] = \frac{f^{(2n+2)}(\cdot)}{(2n+2)!} \int_a^b {}_{n+1}^2(x)(x) dx. \quad (5.8)$$

下面讨论高斯求积公式的稳定性与收敛性.

定理 6 高斯求积公式(5.1)的求积系数 A_k ($k = 0, 1, \dots, n$) 全是正的.

证明 考察

$$l_k(x) = \prod_{\substack{j=0 \\ j \neq k}}^n \frac{x - x_j}{x_k - x_j},$$

它是 n 次多项式, 因而 $l_k(x)$ 是 $2n$ 次多项式, 故高斯求积公式(5.1)对于它能准确成立, 即有

$$0 < \int_a^b l_k^2(x)(x) dx = \sum_{i=0}^n A_i l_k^2(x_i).$$

注意到 $l_k(x_i) = 1$, 上式右端实际上即等于 A_k , 从而有

$$A_k = \int_a^b l_k^2(x)(x) dx > 0.$$

定理得证.

由本定理及定理2, 则得

推论 高斯求积公式(5.1)是稳定的.

定理7 设 $f(x) \in C[a, b]$, 则高斯求积公式(5.1)是收敛的, 即

$$\lim_{n \rightarrow \infty} \sum_{k=0}^n A_k f(x_k) = \int_a^b f(x) dx.$$

证明见[1].

4.5.2 高斯-勒让德求积公式

在高斯求积公式(5.1)中, 若取权函数 $w(x) = 1$, 区间为 $[-1, 1]$, 则得公式

$$\int_{-1}^1 f(x) dx = \sum_{k=0}^n A_k f(x_k). \quad (5.9)$$

我们知道勒让德多项式(见第3章(2.5)式)是区间 $[-1, 1]$ 上的正交多项式, 因此, 勒让德多项式 $P_{n+1}(x)$ 的零点就是求积公式(5.9)的高斯点. 形如(5.9)的高斯公式特别地称为高斯-勒让德求积公式.

若取 $P_1(x) = x$ 的零点 $x_0 = 0$ 做节点构造求积公式

$$\int_{-1}^1 f(x) dx = A_0 f(0).$$

令它对 $f(x) = 1$ 准确成立, 即可定出 $A_0 = 2$. 这样构造出的一点高斯-勒让德求积公式是中矩形公式.

再取 $P_2(x) = \frac{1}{2}(3x^2 - 1)$ 的两个零点 $\pm \frac{1}{\sqrt{3}}$ 构造求积公式

$$\int_{-1}^1 f(x) dx = A_0 f(-\frac{1}{\sqrt{3}}) + A_1 f(\frac{1}{\sqrt{3}}),$$

令它对 $f(x) = 1, x$ 都准确成立, 有

$$A_0 + A_1 = 2;$$

$$A_0 - \frac{1}{3} + A_1 - \frac{1}{3} = 0.$$

由此解出 $A_0 = A_1 = 1$, 从而得到两点高斯-勒让德求积公式

$$\int_{-1}^1 f(x) dx = f(-\frac{1}{3}) + f(\frac{1}{3}).$$

三点高斯-勒让德公式的形式是

$$\int_{-1}^1 f(x) dx = \frac{5}{9}f(-\frac{15}{5}) + \frac{8}{9}f(0) + \frac{5}{9}f(\frac{15}{5}).$$

表 4-7 列出高斯-勒让德求积公式(5.9)的节点和系数。

表 4-7

n	x_k	A_k
0	0.0000000	2.0000000
1	± 0.5773503	1.0000000
2	± 0.7745967	0.5555556
	0.0000000	0.8888889
3	± 0.8611363	0.3478548
	± 0.3399810	0.6521452
4	± 0.9061798	0.2369269
	± 0.5384693	0.4786287
	0.0000000	0.5688889

公式(5.9)的余项由(5.8)得

$$R_n[f] = \frac{\int_{-1}^{2n+2} f(x) dx}{(2n+2)!} L_{n+1}^2(x), \quad [-1, 1],$$

这里 $L_{n+1}(x)$ 是最高项系数为 1 的勒让德多项式, 由第 3 章(2.6)及(2.7)得

$$R_n[f] = \frac{2^{2n+3} [(n+1)!]^4}{(2n+3)[(2n+2)!]^3} f^{(2n+2)}(\) , \quad (-1, 1) . \quad (5.10)$$

当 $n=1$ 时, 有

$$R_1[f] = \frac{1}{135} f^{(4)}(\) .$$

它比辛普森公式余项 $R_1[f] = -\frac{1}{90} f^{(4)}(\)$ (区间为 $[-1, 1]$) 还小,

且比辛普森公式少算一个函数值.

当积分区间不是 $[-1, 1]$, 而是一般的区间 $[a, b]$ 时, 只要做变换

$$x = \frac{b-a}{2}t + \frac{a+b}{2},$$

可将 $[a, b]$ 化为 $[-1, 1]$, 这时

$$\int_a^b f(x) dx = \frac{b-a}{2} \int_{-1}^1 f\left(\frac{b-a}{2}t + \frac{a+b}{2}\right) dt . \quad (5.11)$$

对等式右端的积分即可使用高斯-勒让德求积公式.

例 6 用 4 点 ($n=3$) 的高斯-勒让德求积公式计算

$$I = \int_0^{\frac{\pi}{2}} x^2 \cos x dx .$$

解 先将区间 $[0, \frac{\pi}{2}]$ 化为 $[-1, 1]$, 由 (5.11) 有

$$I = \int_{-1}^1 \frac{1}{4} (1+t)^2 \cos \frac{\pi}{4}(1+t) dt .$$

根据表 4-7 中 $n=3$ 的节点及系数值可求得

$$I = \sum_{k=0}^3 A_k f(x_k) = 0.467402 \quad (\text{准确值 } I = 0.467401 \dots) .$$

4.5.3 高斯-切比雪夫求积公式

若 $a = -1, b = 1$, 且取权函数

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

则所建立的高斯公式为

$$\int_{-1}^1 \frac{f(x)}{1 - x^2} dx = \sum_{k=0}^n A_k f(x_k). \quad (5.12)$$

特别地称为高斯-切比雪夫求积公式。由于区间 $[-1, 1]$ 上关于权函数 $\frac{1}{1 - x^2}$ 的正交多项式是切比雪夫多项式(见第3章第2节)，因此求积公式(5.12)的高斯点是 $n+1$ 次切比雪夫多项式的零点，即为

$$x_k = \cos \frac{2k+1}{2n+2} \quad (k = 0, 1, \dots, n).$$

通过计算(见[2])可知(5.12)的系数 $A_k = \frac{1}{n+1}$ ，使用时将 $n+1$ 个节点公式改为 n 个节点，于是高斯-切比雪夫求积公式写成

$$\int_{-1}^1 \frac{f(x)}{1 - x^2} dx = \sum_{k=1}^n f(x_k), \quad x_k = \cos \frac{(2k-1)}{2n}, \quad (5.13)$$

公式余项由(5.9)可算得

$$R[f] = \frac{2}{2^n (2n)!} f^{(2n)}(\cdot), \quad (-1, 1). \quad (5.14)$$

带权的高斯求积公式可用于计算奇异积分。

例 7 用 5 点($n=5$)的高斯-切比雪夫求积公式计算积分

$$I = \int_{-1}^1 \frac{e^x}{1 - x^2} dx.$$

解 这里 $f(x) = e^x$, $f^{(2n)}(x) = e^x$, 当 $n=5$ 时由公式(5.13)可得

$$I = \frac{1}{5} \sum_{k=1}^5 e^{\cos \frac{2k-1}{10}} = 3.977463.$$

由余项(5.14)可估计误差

$$|R[f]| \leq \frac{e}{2^9 \cdot 10!} \cdot 4.6 \times 10^{-9}.$$

4.6 数值微分

4.6.1 中点方法与误差分析

数值微分就是用函数值的线性组合近似函数在某点的导数值. 按导数定义可以简单地用差商近似导数, 这样立即得到几种数值微分公式

$$\begin{aligned} f'(a) &= \frac{f(a+h) - f(a)}{h}, \\ f'(a) &= \frac{f(a) - f(a-h)}{h}, \\ f'(a) &= \frac{f(a+h) - f(a-h)}{2h}. \end{aligned} \quad (6.1)$$

其中 h 为一增量, 称为步长, 后一种数值微分方法称为中点方法, 它其实是前两种方法的算术平均. 但它的误差阶却由 $O(h)$ 提高到 $O(h^2)$. 上面给出的三个公式是很实用的. 尤其是中点公式更为常用.

为要利用中点公式

$$G(h) = \frac{f(a+h) - f(a-h)}{2h}$$

计算导数 $f'(a)$ 的近似值, 首先必须选取合适的步长, 为此需要进行误差分析. 分别将 $f(a \pm h)$ 在 $x=a$ 处做泰勒展开有

$$\begin{aligned} f(a \pm h) &= f(a) \pm hf'(a) + \frac{h^2}{2!}f''(a) \pm \frac{h^3}{3!}f'''(a) \\ &\quad + \frac{h^4}{4!}f^{(4)}(a) \pm \frac{h^5}{5!}f^{(5)}(a) + \dots \end{aligned}$$

代入上式得

$$G(h) = f(a) + \frac{h^2}{3!} f'(a) + \frac{h^4}{5!} f^{(5)}(a) + \dots$$

由此得知, 从截断误差的角度看, 步长越小, 计算结果越准确. 且

$$|f(a) - G(h)| \leq \frac{h^2}{6} M, \quad (6.2)$$

其中 $M = \max_{|x-a|} |f(x)|$.

再考察舍入误差. 按中点公式计算, 当 h 很小时, 因 $f(a+h)$ 与 $f(a-h)$ 很接近, 直接相减会造成有效数字的严重损失 (参看第 1 章第 4 节). 因此, 从舍入误差的角度来看, 步长是不宜太小的.

例如, 用中点公式求 $f(x) = x$ 在 $x=2$ 处的一阶导数

$$G(h) = \frac{2+h-2-h}{2h}.$$

设取 4 位数字计算. 结果见表 4-8 (导数的准确值 $f(2) = 0.353553$).

表 4-8

h	$G(h)$	h	$G(h)$	h	$G(h)$
1	0.3660	0.05	0.3530	0.001	0.3500
0.5	0.3564	0.01	0.3500	0.0005	0.3000
0.1	0.3535	0.005	0.3500	0.0001	0.3000

从表 4-8 中看到 $h=0.1$ 的逼近效果最好, 如果进一步缩小步长, 则逼近效果反而越差. 这里因为当 $f(a+h)$ 及 $f(a-h)$ 分别有差入误差 $|_1$ 及 $|_2$. 若令 $= \max\{|_1|, |_2|\}$, 则计算 $f(a)$ 的舍入误差上界为

$$(f(a)) = |f(a) - G(a)| \leq \frac{|_1| + |_2|}{2h} = \frac{h}{h},$$

它表明 h 越小, 舍入误差 $(f(a))$ 越大, 故它是病态的. 用中点公式(6.1)计算 $f(a)$ 的误差上界为

$$E(h) = \frac{h^2}{6} M + \frac{1}{h},$$

要使误差 $E(h)$ 最小, 步长 h 不宜太大, 也不宜太小. 其最优步长应为 $h_{\text{opt}} = \sqrt[3]{3/M}$.

4.6.2 插值型的求导公式

对于列表函数 $y = f(x)$:

x	x_0	x_1	x_2	...	x_n
y	y_0	y_1	y_2	...	y_n

运用插值原理, 可以建立插值多项式 $y = P_n(x)$ 作为它的近似. 由于多项式的求导比较容易, 我们取 $P_n(x)$ 的值作为 $f(x)$ 的近似值, 这样建立的数值公式

$$f(x) = P_n(x) \quad (6.3)$$

统称插值型的求导公式.

必须指出, 即使 $f(x)$ 与 $P_n(x)$ 的值相差不多, 导数的近似值 $P_n'(x)$ 与导数的真值 $f'(x)$ 仍然可能差别很大, 因而在使用求导公式(6.3)时应特别注意误差的分析.

依据插值余项定理, 求导公式(6.3)的余项为

$$\begin{aligned} f(x) - P_n(x) &= \frac{f^{(n+1)}(\cdot)}{(n+1)!} {}_{n+1}(x) \\ &\quad + \frac{f^{(n+1)}(\cdot)}{(n+1)!} \frac{d}{dx} f^{(n+1)}(\cdot), \end{aligned}$$

式中

$${}_{n+1}(x) = \prod_{i=0}^n (x - x_i).$$

在这一余项公式中, 由于 x 是 x 的未知函数, 我们无法对它的第二项 $\frac{n+1}{(n+1)!} \frac{d}{dx} f^{(n+1)}(\cdot)$ 做出进一步的说明。因此, 对于随意给出的点 x , 误差 $f(x) - P_n(x)$ 是无法预估的。但是, 如果我们限定求某个节点 x_k 上的导数值, 那么上面的第二项因式 $\frac{n+1}{(n+1)!} f^{(n+1)}(x_k)$ 变为零, 这时有余项公式

$$f(x_k) - P_n(x_k) = \frac{f^{(n+1)}(\cdot)}{(n+1)!} |_{n+1}(x_k). \quad (6.4)$$

下面我们仅仅考察节点处的导数值。为简化讨论, 假定所给的节点是等距的。

1. 两点公式

设已给出两个节点 x_0, x_1 上的函数值 $f(x_0), f(x_1)$, 做线性插值得公式

$$P_1(x) = \frac{x - x_1}{x_0 - x_1} f(x_0) + \frac{x - x_0}{x_1 - x_0} f(x_1).$$

对上式两端求导, 记 $x_1 - x_0 = h$, 有

$$P_1'(x) = \frac{1}{h} [-f(x_0) + f(x_1)],$$

于是有下列求导公式:

$$P_1'(x_0) = \frac{1}{h} [f(x_1) - f(x_0)];$$

$$P_1'(x_1) = \frac{1}{h} [f(x_1) - f(x_0)].$$

而利用余项公式(6.4)知, 带余项的两点公式是

$$f(x_0) = \frac{1}{h} [f(x_1) - f(x_0)] - \frac{h}{2} f'(\cdot);$$

$$f(x_1) = \frac{1}{h} [f(x_1) - f(x_0)] + \frac{h}{2} f'(\cdot).$$

2. 三点公式

设已给出三个节点 $x_0, x_1 = x_0 + h, x_2 = x_0 + 2h$ 上的函数值,

做二次插值

$$\begin{aligned} P_2(x) &= \frac{(x - x_0)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} f(x_0) \\ &\quad + \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)} f(x_1) \\ &\quad + \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)} f(x_2), \end{aligned}$$

令 $x = x_0 + th$, 上式可表示为

$$\begin{aligned} P_2(x_0 + th) &= \frac{1}{2}(t-1)(t-2)f(x_0) - t(t-2)f(x_1) \\ &\quad + \frac{1}{2}t(t-1)f(x_2). \end{aligned}$$

两端对 t 求导, 有

$$\begin{aligned} P_2(x_0 + th) &= \frac{1}{2h}[(2t-3)f(x_0) - (4t-4)f(x_1) \\ &\quad + (2t-1)f(x_2)]. \end{aligned} \quad (6.5)$$

这里撇号()表示对变量 x 求导数. 上式分别取 $t=0, 1, 2$, 得到三种三点公式:

$$\begin{aligned} P_2(x_0) &= \frac{1}{2h}[-3f(x_0) + 4f(x_1) - f(x_2)]; \\ P_2(x_1) &= \frac{1}{2h}[-f(x_0) + f(x_2)]; \\ P_2(x_2) &= \frac{1}{2h}[f(x_0) - 4f(x_1) + 3f(x_2)]. \end{aligned}$$

而带余项的三点求导公式如下:

$$\begin{aligned} f(x_0) &= \frac{1}{2h}[-3f(x_0) + 4f(x_1) - f(x_2)] + \frac{h^2}{3}f''(); \\ f(x_1) &= \frac{1}{2h}[-f(x_0) + f(x_2)] - \frac{h^2}{6}f''(); \\ f(x_2) &= \frac{1}{2h}[f(x_0) - 4f(x_1) + 3f(x_2)] + \frac{h^2}{3}f''(). \end{aligned} \quad (6.6)$$

其中的公式(6.6)是我们所熟悉的中点公式. 在三点公式中, 它由于少用了一个函数值 $f(x_1)$ 而引人注目.

用插值多项式 $P_n(x)$ 作为 $f(x)$ 的近似函数, 还可以建立高阶数值微分公式:

$$f^{(k)}(x) = P_m^{(k)}(x), \quad k = 1, 2, \dots$$

例如, 将式(6.5)再对 t 求导一次, 有

$$P_2(x_0 + th) = \frac{1}{h^2} [f(x_0) - 2f(x_1) + f(x_2)],$$

于是有

$$P_2(x_1) = \frac{1}{h^2} [f(x_1 - h) - 2f(x_1) + f(x_1 + h)].$$

而带余项的二阶三点公式如下:

$$f(x_1) = \frac{1}{h^2} [f(x_1 - h) - 2f(x_1) + f(x_1 + h)] - \frac{h^2}{12} f^{(4)}() . \quad (6.7)$$

4.6.3 利用数值积分求导

微分是积分的逆运算, 因此可利用数值积分的方法来计算数值微分. 设 $f(x)$ 是一个充分光滑的函数, 设 $(x) = f(x)$, $x_k = a + kh$, $k = 0, 1, \dots, n$, $h = \frac{b-a}{n}$, 则有

$$f(x_{k+1}) = f(x_{k-1}) + \int_{x_{k-1}}^{x_{k+1}} (x) dx \quad (k = 1, \dots, n-1), \quad (6.8)$$

对上式右边积分采用不同的求积公式就可得到不同的数值微分公式. 例如, 对 $\int_{x_{k-1}}^{x_{k+1}} (x) dx$ 用中矩形公式(1.2), 则得

$$\int_{x_{k-1}}^{x_{k+1}} (x) dx = 2h (x_k) + \frac{1}{24} (2h)^3 () , \quad k = (x_{k-1}, x_{k+1}).$$

从而得到中点微分公式

$$f(x_k) = \frac{f(x_{k+1}) - f(x_{k-1})}{2h} - \frac{1}{6}h^2 f''(x_k).$$

若对(6.8)右端积分用辛普森求积公式, 则有

$$\int_{x_{k-1}}^{x_{k+1}} f(x) dx = \frac{h}{3} [f(x_{k-1}) + 4f(x_k) + f(x_{k+1})] - \frac{h^5}{90} f^{(4)}(x_k), \quad k \in (x_{k-1}, x_{k+1}),$$

上式略去余项, 并记 $f(x_k)$ 的近似值为 m_k , 则得到辛普森数值微分公式

$$m_{k-1} + 4m_k + m_{k+1} = \frac{3}{h} [f(x_{k+1}) - f(x_{k-1})] \quad (k = 1, \dots, n-1).$$

这是关于 m_0, m_1, \dots, m_n 这 $n+1$ 个未知量的 $n-1$ 个方程组, 若 $m_0 = f(x_0), m_n = f(x_n)$ 已知, 则可得

$$\begin{array}{ccccccccc} 4 & 1 & & m_1 & & \frac{3}{h} [f(x_2) - f(x_0)] & - & f(x_0) \\ & & & & & & & \\ 1 & 4 & 1 & m_2 & & \frac{3}{h} [f(x_3) - f(x_1)] & & \\ & & & & & & & \\ w & w & w & \dots & = & \dots & & \\ & & & & & & & \\ 1 & 4 & 1 & m_{n-2} & & \frac{3}{h} [f(x_{n-1}) - f(x_{n-3})] & & \\ & & & & & & & \\ 1 & 4 & m_{n-1} & & & \frac{3}{h} [f(x_n) - f(x_{n-2})] & - & f(x_n) \end{array} \quad (6.9)$$

这是关于 m_1, \dots, m_{n-1} 的三对角方程组, 且系数矩阵为严格对角占优的, 可用追赶法求解(见第5章5.4节).

如果端点导数值不知道, 那么对(6.9)中第1个和第 $n-1$ 个方程可分别用 $f(x_1)$ 及 $f(x_{n-1})$ 的中点微分公式近似, 即取

$$m_1 = \frac{1}{2h} [f(x_2) - f(x_0)], \quad m_{n-1} = \frac{1}{2h} [f(x_n) - f(x_0)].$$

然后求 m_2, \dots, m_{n-2} 即为 $f(x_2), \dots, f(x_{n-2})$ 的近似值.

例 8 给定 $f(x) = \sqrt{x}$ 的一张数据表(表 4-9 左部), 并给定 $f(100)$ 及 $f(105)$ 的值(见表 4-9). 利用辛普森数值微分公式求 $f'(x)$ 在 $x = 101, 102, 103, 104$ 上的一阶导数.

解 根据(6.9)有

$$\begin{array}{cccc} 4 & 1 & m_1 & 0.24851482 \\ 1 & 4 & 1 & m_2 = 0.29704785 \\ & 1 & 4 & m_3 = 0.29560227 \\ & 1 & 4 & m_4 = 0.24538260 \end{array} .$$

解之得 m_i ($i = 1, 2, 3, 4$), 结果见表 4-9.

表 4-9

k	x_k	$f(x_k) = x_k$	$f(x_k)$	$m_k f(x_k)$
0	100	10.00000000	0.05000000	
1	101	10.04987562	0.049751859	0.04975186
2	102	10.09950494	0.049507377	0.049507376
3	103	10.14889157	0.049266463	0.049266463
4	104	10.19803903	0.049029033	0.049029033
5	105	10.24695077	0.048795003	

4.6.4 三次样条求导

三次样条函数 $S(x)$ 作为 $f(x)$ 的近似, 不但函数值很接近, 导数值也很接近, 并有

$$f^{(k)}(x) - S^{(k)}(x) = C_k f^{(4)}(x) h^{4-k} \quad (k = 0, 1, 2). \quad (6.10)$$

(见第 2 章定理 4), 因此利用三次样条函数 $S(x)$ 直接得到

$$f^{(k)}(x) = S^{(k)}(x) \quad (k = 0, 1, 2).$$

根据第 2 章(7.8), (7.9)可求得

$$f(x_k) - S(x_k) = -\frac{h_k}{3} M_k - \frac{h_k}{6} M_{k+1} + f[x_k, x_{k+1}],$$

$$f(x_k) = M_k.$$

这里 $f[x_k, x_{k+1}]$ 为一阶均差. 其误差由(6.10)可得

$$f - S \quad \frac{1}{24} f^{(4)} h^3,$$

$$f - S \quad \frac{3}{8} f^{(4)} h^2.$$

4.6.5 数值微分的外推算法

利用中点公式计算导数值时

$$f(x) - G(h) = \frac{1}{2h} [f(x+h) - f(x-h)].$$

对 $f(x)$ 在点 x 做泰勒级数展开有

$$f(x) = G(h) + {}_1 h^2 + {}_2 h^4 + \dots,$$

其中 ${}_i$ ($i=1, 2, \dots$) 与 h 无关, 利用理查森外推(见本章第4节)对 h 逐次分半, 若记 $G_0(h) = G(h)$, 则有

$$G_m(h) = \frac{4^m G_{m-1} - \frac{h}{2} G_{m-1}(h)}{4^m - 1} \quad (m=1, 2, \dots). \quad (6.11)$$

公式(6.11)的计算过程见表 4-10, 表中 i 为外推步数.

表 4-10

$$G(h)$$

$$G \quad \frac{h}{2} \quad G_1(h)$$

$$G \quad \frac{h}{2^2} \quad G_1 \quad \frac{h}{2} \quad G_2(h)$$

$$G \quad \frac{h}{2^3} \quad G_1 \quad \frac{h}{2^2} \quad G_2 \quad \frac{h}{2} \quad G_3(h)$$

...

...

...

...

W

根据理查森外推方法, (6.11)的误差为

$$f(x) - G_m(h) = O(h^{2(m+1)}).$$

由此看出当 m 较大时, 计算是很精确的. 考虑到舍入误差, 一般 m 不能取太大.

例 9 用外推法计算 $f(x) = x^2 e^{-x}$ 在 $x=0.5$ 的导数.

$$\text{解} \quad \text{令 } G(h) = \frac{1}{2h} \left(\frac{1}{2} + h \right)^2 e^{-\frac{1}{2}+h} - \frac{1}{2} \left(h \right)^2 e^{-\frac{1}{2}-h},$$

当 $h=0.1, 0.05, 0.025$ 时, 由外推法表 4-10 可算得

$$G(0.1) = 0.4516049081$$

$$G(0.05) = 0.4540761693 \quad G(h) = 0.4548999231$$

$$G(0.025) = 0.4546926288 \quad G \frac{h}{2} = 0.4548981152$$

$$G = 0.454897994$$

$f(0.5)$ 的精确值为 0.454897994, 可见当 $h=0.025$ 时用中点微分公式只有 3 位有效数字, 外推一次达到 5 位有效数字, 外推两次达到 9 位有效数字.

评注

本章介绍积分和微分的数值计算方法. 我们知道, 积分和微分是两种分析运算, 它们都是用极限来定义的. 数值积分和数值微分则归结为函数值的四则运算, 从而使计算过程可以在计算机上完成.

处理数值积分和数值微分的基本方法是逼近法: 设法构造某个简单函数 $P(x)$ 近似 $f(x)$, 然后对 $P(x)$ 求积(求导)得到 $f(x)$ 的积分(导数)的近似值. 本章基于插值原理推导了数值积分和数值微分的基本公式.

插值求积公式分牛顿-柯特斯公式和高斯公式两类. 前者为

等距节点,但 $n=8$ 时计算不稳定,实际计算宜采用复合求积方法. 高斯公式精度高,计算稳定,但节点选取较困难. 带权高斯求积方法,能把复杂积分化简,还可以直接计算奇异积分. 这两类公式也可通过求积公式的代数精度建立.

基于理查森外推的龙贝格求积方法由于计算程序简单,精度较高,是一个在计算机上求积的有效算法. 在数值微分中也有相似的算法,外推方法和思想是数值分析中一种很重要的方法.

数值微分由于计算不稳定性,步长选取是很重要的. 通常数值微分的外推法可得到较满意结果,但 h 也不能太小.

数值积分中一些重要内容如奇异积分,振荡函数积分和二重积分计算等均未涉及,可参见文献[1].

习题

1. 确定下列求积公式中的待定参数,使其代数精度尽量高,并指明所构造出的求积公式所具有的代数精度:

$$1) \int_{-h}^h f(x) dx = A_{-1} f(-h) + A_0 f(0) + A_1 f(h);$$

$$2) \int_{-2h}^{2h} f(x) dx = A_{-1} f(-h) + A_0 f(0) + A_1 f(h);$$

$$3) \int_{-1}^1 f(x) dx = [f(-1) + 2f(x_1) + 3f(x_2)]/3;$$

$$4) \int_0^h f(x) dx = h[f(0) + f(h)]/2 + ah^2[f(0) - f(h)].$$

2. 分别用梯形公式和辛普森公式计算下列积分:

$$1) \int_0^1 \frac{x}{4+x^2} dx, \quad n=8;$$

$$2) \int_0^1 \frac{(1-e^{-x})^{\frac{1}{2}}}{x} dx, \quad n=10;$$

$$3) \int_1^9 x dx, \quad n=4;$$

4) $\int_0^6 4 - \sin^2 d \, dx, \quad n = 6.$

3. 直接验证柯特斯公式(2.4)具有5次代数精度.

4. 用辛普森公式求积分 $\int_0^1 e^{-x} dx$ 并估计误差.

5. 推导下列三种矩形求积公式:

$$\int_a^b f(x) dx = (b - a) f(a) + \frac{f(\frac{a+b}{2})}{2} (b - a)^2;$$

$$\int_a^b f(x) dx = (b - a) f(b) - \frac{f(\frac{a+b}{2})}{2} (b - a)^2;$$

$$\int_a^b f(x) dx = (b - a) f\left(\frac{a+b}{2}\right) + \frac{f(\frac{a+b}{4})}{24} (b - a)^3.$$

6. 若用复化梯形公式计算积分 $I = \int_0^1 e^x dx$, 问区间 $[0, 1]$ 应分多少等分

才能使截断误差不超过 $\frac{1}{2} \times 10^{-5}$? 若改用复化辛普森公式, 要达到同样精度区间 $[0, 1]$ 应分多少等分?

7. 如果 $f(x) > 0$, 证明用梯形公式计算积分 $I = \int_a^b f(x) dx$ 所得结果比准确值 I 大, 并说明其几何意义.

8. 用龙贝格求积方法计算下列积分, 使误差不超过 10^{-5} .

(1) $\int_0^1 e^{-x} dx,$

(2) $\int_0^2 x \sin x dx,$

(3) $\int_0^3 x - 1 + x^2 dx.$

9. 用 $n=2, 3$ 的高斯-勒让德公式计算积分

$$\int_1^3 e^x \sin x dx.$$

10. 地球卫星轨道是一个椭圆, 椭圆周长的计算公式是

$$S = a \int_0^{2\pi} 1 - \frac{c}{a} \sin^2 d \, d\theta,$$

这里 a 是椭圆的半长轴, c 是地球中心与轨道中心(椭圆中心)的距离, 记 h 为近地点距离, H 为远地点距离, $R=6371(\text{km})$ 为地球半径, 则

$$a = (2R + H + h)/2, \quad c = (H - h)/2.$$

我国第一颗人造地球卫星近地点距离 $h = 439$ (km), 远地点距离 $H = 2384$ (km), 试求卫星轨道的周长 .

11. 证明等式

$$n \sin \frac{\pi}{n} = 1 - \frac{3}{3!n^2} + \frac{5}{5!n^4} - \dots$$

试依据 $n \sin(\pi/n)$ ($n=3, 6, 12$) 的值, 用外推算法求 π 的近似值 .

12. 用下列方法计算积分 $\int_1^3 \frac{dy}{y}$, 并比较结果 .

1) 龙贝格方法;

2) 三点及五点高斯公式;

3) 将积分区间分为四等分, 用复化两点高斯公式 .

13. 用三点公式和积分方法求 $f(x) = \frac{1}{(1+x)^2}$ 在 $x = 1.0, 1.1$ 和 1.2 处

的导数值, 并估计误差 . $f(x)$ 的值由下表给出 :

x	1.0	1.1	1.2
$f(x)$	0.2500	0.2268	0.2066

第 5 章 解线性方程组的直接方法

5.1 引言与预备知识

5.1.1 引言

在自然科学和工程技术中很多问题的解决常常归结为解线性代数方程组,例如电学中的网络问题,船体数学放样中建立三次样条函数问题,用最小二乘法求实验数据的曲线拟合问题,解非线性方程组问题,用差分法或者有限元方法解常微分方程、偏微分方程边值问题等都导致求解线性代数方程组,而这些方程组的系数矩阵大致分为两种,一种是低阶稠密矩阵(例如,阶数不超过 150),另一种是大型稀疏矩阵(即矩阵阶数高且零元素较多)。

关于线性方程组的数值解法一般有两类:

1. 直接法

就是经过有限步算术运算,可求得方程组精确解的方法(若计算过程中没有舍入误差)。但实际计算中由于舍入误差的存在和影响,这种方法也只能求得线性方程组的近似解。本章将阐述这类算法中最基本的高斯消去法及其某些变形。这类方法是解低阶稠密矩阵方程组及某些大型稀疏方程组(例如,大型带状方程组)的有效方法。

2. 迭代法

就是用某种极限过程去逐步逼近线性方程组精确解的方法。迭代法具有需要计算机的存贮单元较少、程序设计简单、原始系数矩阵在计算过程中始终不变等优点,但存在收敛性及收敛速度问

题. 迭代法是解大型稀疏矩阵方程组(尤其是由微分方程离散后得到的大型方程组)的重要方法(见第 6 章).

为了讨论线性方程组数值解法, 需复习一些基本的矩阵代数知识.

5.1.2 向量和矩阵

用 $\mathbf{R}^{m \times n}$ 表示全部 $m \times n$ 实矩阵的向量空间, $\mathbf{C}^{m \times n}$ 表示全部 $m \times n$ 复矩阵的向量空间.

$$\mathbf{A} \in \mathbf{R}^{m \times n} \quad \mathbf{A} = (a_{ij}) = \begin{array}{cccc} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{array}$$

(实数排成的矩形表, 称为 m 行 n 列矩阵).

$$\mathbf{x} \in \mathbf{R}^n \quad \mathbf{x} = \begin{array}{c} x_1 \\ x_2 \\ \dots \\ x_n \end{array} \quad (n \text{ 维列向量}).$$

$$\mathbf{A} = (\mathbf{a}_1 \quad \mathbf{a}_2 \quad \dots \quad \mathbf{a}_n),$$

其中 \mathbf{a}_i 为 \mathbf{A} 的第 i 列. 同理

$$\mathbf{b}^T$$

$$\mathbf{A} = \begin{array}{c} \dots \\ \mathbf{b}_1^T \\ \mathbf{b}_2^T \\ \dots \\ \mathbf{b}_m^T \end{array},$$

其中 \mathbf{b}_i^T 为 \mathbf{A} 的第 i 行.

矩阵的基本运算:

(1) 矩阵加法 $\mathbf{C} = \mathbf{A} + \mathbf{B}$, $c_{ij} = a_{ij} + b_{ij}$ ($\mathbf{A} \in \mathbf{R}^{m \times n}$, $\mathbf{B} \in \mathbf{R}^{m \times n}$, $\mathbf{C} \in \mathbf{R}^{m \times n}$).

(2) 矩阵与标量的乘法 $\mathbf{C} = k\mathbf{A}$, $c_{ij} = ka_{ij}$.

(3) 矩阵与矩阵乘法 $\mathbf{C} = \mathbf{AB}$, $c_{ij} = \sum_{k=1}^n a_{ik} b_{kj}$ ($\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{n \times p}$, $\mathbf{C} \in \mathbb{R}^{m \times p}$).

(4) 转置矩阵 $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{C} = \mathbf{A}^T$, $c_{ij} = a_{ji}$.

(5) 单位矩阵 $\mathbf{I} = (\mathbf{e} \ \mathbf{e} \dots \mathbf{e}) \in \mathbb{R}^{n \times n}$, 其中

$$\mathbf{e}_k = (0, \dots, 0, 1, 0, \dots, 0)^T, \quad k = 1, 2, \dots, n.$$

(6) 非奇异矩阵 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\mathbf{B} \in \mathbb{R}^{n \times n}$. 如果 $\mathbf{AB} = \mathbf{BA} = \mathbf{I}$, 则称 \mathbf{B} 是 \mathbf{A} 的逆矩阵, 记为 \mathbf{A}^{-1} , 且 $(\mathbf{A}^{-1})^T = (\mathbf{A}^T)^{-1}$. 如果 \mathbf{A}^{-1} 存在, 则称 \mathbf{A} 为非奇异矩阵. 如果 $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$ 均为非奇异矩阵, 则 $(\mathbf{AB})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1}$.

(7) 矩阵的行列式 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$, 则 \mathbf{A} 的行列式可按任一行(或列)展开, 即

$$\det(\mathbf{A}) = \sum_{j=1}^n a_{ij} A_{ij} \quad (i = 1, 2, \dots, n),$$

其中 A_{ij} 为 a_{ij} 的代数余子式, $A_{ij} = (-1)^{i+j} M_{ij}$, M_{ij} 为元素 a_{ij} 的余子式.

行列式性质:

- (a) $\det(\mathbf{AB}) = \det(\mathbf{A}) \det(\mathbf{B})$, $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$.
- (b) $\det(\mathbf{A}^T) = \det(\mathbf{A})$, $\mathbf{A} \in \mathbb{R}^{n \times n}$.
- (c) $\det(c\mathbf{A}) = c^n \det(\mathbf{A})$, $c \in \mathbb{R}$, $\mathbf{A} \in \mathbb{R}^{n \times n}$.
- (d) $\det(\mathbf{A}) = 0$ \mathbf{A} 是非奇异矩阵.

5.1.3 特殊矩阵

设 $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{n \times n}$.

- (1) 对角矩阵 如果当 $i = j$ 时, $a_{ij} = 0$.
- (2) 三对角矩阵 如果当 $|i - j| > 1$ 时, $a_{ij} = 0$.
- (3) 上三角矩阵 如果当 $i > j$ 时, $a_{ij} = 0$.
- (4) 上海森伯格(Hessenberg)阵 如果当 $i > j + 1$ 时, a_{ij}

$= 0$.

(5) 对称矩阵 如果 $\mathbf{A}^T = \mathbf{A}$.

(6) 埃尔米特矩阵 设 $\mathbf{A} \in \mathbb{C}^{n \times n}$, 如果 $\mathbf{A}^H = \mathbf{A}$ ($\mathbf{A}^H = \overline{\mathbf{A}}^T$, 即为 \mathbf{A} 的共轭转置).

(7) 对称正定矩阵 如果 (a) $\mathbf{A}^T = \mathbf{A}$, (b) 对任意非零向量 $\mathbf{x} \in \mathbb{R}^n$, $(\mathbf{A}\mathbf{x}, \mathbf{x}) = \mathbf{x}^T \mathbf{A}\mathbf{x} > 0$.

(8) 正交矩阵 如果 $\mathbf{A}^{-1} = \mathbf{A}^T$.

(9)酉矩阵 设 $\mathbf{A} \in \mathbb{C}^{n \times n}$, 如果 $\mathbf{A}^{-1} = \mathbf{A}^H$.

(10) 初等置换阵 由单位矩阵 \mathbf{I} 交换第 i 行与第 j 行(或交换第 i 列与第 j 列), 得到的矩阵记为 \mathbf{I}_{ij} , 且

$\mathbf{I}_{ij}\mathbf{A} = \mathbf{A}'$ (为交换 \mathbf{A} 第 i 行与第 j 行得到的矩阵);

$\mathbf{A}\mathbf{I}_j = \mathbf{B}$ (为交换 \mathbf{A} 第 i 列与第 j 列得到的矩阵).

(11) 置换阵 由初等置换阵的乘积得到的矩阵.

定理 1 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$, 则下述命题等价:

(1) 对任何 $\mathbf{b} \in \mathbb{R}^n$, 方程组 $\mathbf{A}\mathbf{x} = \mathbf{b}$ 有唯一解.

(2) 齐次方程组 $\mathbf{A}\mathbf{x} = \mathbf{0}$ 只有唯一解 $\mathbf{x} = \mathbf{0}$.

(3) $\det(\mathbf{A}) \neq 0$.

(4) \mathbf{A}^{-1} 存在.

(5) \mathbf{A} 的秩 $\text{rank}(\mathbf{A}) = n$.

定理 2 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 为对称正定阵, 则

(1) \mathbf{A} 为非奇异矩阵, 且 \mathbf{A}^{-1} 亦是对称正定阵.

(2) 记 \mathbf{A}_k 为 \mathbf{A} 的顺序主子阵, 则 \mathbf{A}_k ($k = 1, 2, \dots, n$) 亦是对称正定矩阵, 其中

$$\mathbf{A}_k = \begin{array}{ccc|c} a_{11} & \cdots & a_{1k} \\ \vdots & \ddots & \vdots & (k=1, 2, \dots, n) \\ a_{kk} & \cdots & a_{kk} \end{array}$$

(3) \mathbf{A} 的特征值 $\lambda_i > 0$ ($i = 1, 2, \dots, n$).

(4) \mathbf{A} 的顺序主子式都大于零, 即 $\det(\mathbf{A}_k) > 0$ ($k = 1, 2, \dots, n$).

定理3 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 为对称矩阵. 如果 $\det(\mathbf{A}) > 0$ ($k = 1, 2, \dots, n$), 或 \mathbf{A} 的特征值 $\lambda_i > 0$ ($i = 1, 2, \dots, n$), 则 \mathbf{A} 为对称正定阵.

有重特征值的矩阵不一定相似于对角矩阵, 那么一般 n 阶矩阵 \mathbf{A} 在相似变换下能简化到什么形状.

定理4(Jordan 标准型) 设 \mathbf{A} 为 n 阶矩阵, 则存在一个非奇异矩阵 \mathbf{P} 使得

$$\mathbf{P}^{-1} \mathbf{A} \mathbf{P} = \begin{matrix} \mathbf{J}_1(1) \\ \mathbf{J}_2(2) \\ \vdots \\ \mathbf{J}_r(r) \end{matrix},$$

其中

$$\begin{aligned} \mathbf{J}_i &= \begin{matrix} i & & & & \\ & \mathbf{W} & & & \\ & & \mathbf{W} & & \\ & & & \ddots & \\ & & & & i \end{matrix}, \\ n_i &= 1 (i = 1, 2, \dots, r), \text{ 且 } \sum_{i=1}^r n_i = n. \end{aligned}$$

为若当(Jordan)块.

(1) 当 \mathbf{A} 的若当标准型中所有若当块 \mathbf{J}_i 均为一阶时, 此标准型变成对角矩阵.

(2) 如果 \mathbf{A} 的特征值各不相同, 则其若当标准型必为对角阵 $\text{diag}(1, 2, \dots, n)$.

5.2 高斯消去法

本节介绍高斯消去法(逐次消去法)及消去法和矩阵三角分解

之间的关系. 虽然高斯消去法是一个古老的求解线性方程组的方法(早在公元前 250 年我国就掌握了解方程组的消去法), 但由它改进、变形得到的选主元素消去法、三角分解法仍然是目前计算机上常用的有效方法.

5.2.1 高斯消去法

设有线性方程组

$$\begin{aligned} a_{11} x_1 + a_{12} x_2 + \dots + a_{1n} x_n &= b_1, \\ a_{21} x_1 + a_{22} x_2 + \dots + a_{2n} x_n &= b_2, \\ &\dots \\ a_{m1} x_1 + a_{m2} x_2 + \dots + a_{mn} x_n &= b_m. \end{aligned} \quad (2.1)$$

或写为矩阵形式

$$\begin{array}{cccccc} a_{11} & a_{12} & \dots & a_{1n} & x_1 & b \\ a_{21} & a_{22} & \dots & a_{2n} & x_2 & = b \\ \dots & \dots & & \dots & \dots & \dots \\ a_{m1} & a_{m2} & \dots & a_{mn} & x_n & b_m \end{array}$$

简记为 $\mathbf{Ax} = \mathbf{b}$.

首先举一个简单的例子来说明消去法的基本思想.

例 1 用消去法解方程组

$$x_1 + x_2 + x_3 = 6, \quad (2.2)$$

$$4x_2 - x_3 = 5, \quad (2.3)$$

$$2x_1 - 2x_2 + x_3 = 1. \quad (2.4)$$

解 第 1 步. 将方程(2.2)乘上 -2 加到方程(2.4)上去, 消去(2.4)中的未知数 x_1 , 得到

$$-4x_2 - x_3 = -11. \quad (2.5)$$

第 2 步. 将方程(2.3)加到方程(2.5)上去, 消去方程(2.5)中的未知数 x_2 , 得到与原方程组等价的三角形方程组

$$\begin{aligned}
 x_1 + x_2 + x_3 &= 6, \\
 4x_2 - x_3 &= 5, \\
 -2x_3 &= -6.
 \end{aligned} \tag{2.6}$$

显然, 方程组(2.6)是容易求解的, 解为

$$x^* = (1, 2, 3)^T.$$

上述过程相当于

$$\begin{array}{ccc|ccc|c}
 1 & 1 & 1 & 6 & 1 & 1 & 1 & 6 \\
 (\mathbf{A}/\mathbf{b}) = & 0 & 4 & -1 & 5 & 0 & 4 & -1 \\
 & 2 & -2 & 1 & 1 & 0 & -4 & -1 & -11 \\
 & 1 & 1 & 1 & 6 & & & \\
 & 0 & 4 & -1 & 5 & & & \\
 & 0 & 0 & -2 & -6 & & & \\
 (-2) \times r_1 + r_3 & & & & r_3 & r_2 + r_3 & r_3
 \end{array}$$

其中用 r_i 表示矩阵的第 i 行.

由此看出, 用消去法解方程组的基本思想是用逐次消去未知数的方法把原方程组 $\mathbf{Ax} = \mathbf{b}$ 化为与其等价的三角形方程组, 而求解三角形方程组可用回代的方法求解. 换句话说, 上述过程就是用行的初等变换将原方程组系数矩阵化为简单形式(上三角矩阵), 从而将求解原方程组(2.1)的问题转化为求解简单方程组的问题. 或者说, 对系数矩阵 \mathbf{A} 施行一些左变换(用一些简单矩阵)将其约化为上三角矩阵.

下面我们讨论求解一般线性方程组的高斯消去法.

将(2.1)记为 $\mathbf{A}^{(1)} \mathbf{x} = \mathbf{b}^{(1)}$, 其中

$$\mathbf{A}^{(1)} = (a_j^{(1)}) = (a_{ij}), \quad \mathbf{b}^{(1)} = \mathbf{b}.$$

(1) 第 1 步($k=1$).

设 $a_1^{(1)} \neq 0$, 首先计算乘数

$$m_{ii} = a_{ii}^{(1)} / a_{11}^{(1)} \quad (i = 2, 3, \dots, m).$$

用 $-m_{ii}$ 乘(2.1)的第一个方程, 加到第 i 个($i = 2, 3, \dots, m$)方程

上, 消去(2.1)的从第二个方程到第 m 个方程中的未知数 x_1 , 得到与(2.1)等价的方程组

$$\begin{array}{cccccc} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1n}^{(1)} & x_1 & b_1^{(1)} \\ 0 & a_{22}^{(2)} & \dots & a_{2n}^{(2)} & x_2 & b_2^{(2)} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & a_{m2}^{(2)} & \dots & a_{mn}^{(2)} & x_n & b_m^{(2)} \end{array} = . \quad (2.7)$$

简记为

$$\mathbf{A}^{(2)} \mathbf{x} = \mathbf{b}^{(2)},$$

其中 $\mathbf{A}^{(2)}$, $\mathbf{b}^{(2)}$ 的元素计算公式为

$$\begin{aligned} a_{ij}^{(2)} &= a_{ij}^{(1)} - m_{i1} a_{1j}^{(1)} \quad (i = 2, \dots, m; j = 2, \dots, n), \\ b_i^{(2)} &= b_i^{(1)} - m_{i1} b_1^{(1)} \quad (i = 2, \dots, m). \end{aligned}$$

(2) 第 k 次消元($k = 1, 2, \dots, s = \min(m - 1, n)$)。

设上述第 1 步, ..., 第 $k - 1$ 步消元过程计算已经完成, 即已计算好与(2.1)等价的方程组

$$\begin{array}{cccccc} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1k}^{(1)} & \dots & a_{1n}^{(1)} & x_1 & b_1^{(1)} \\ a_{22}^{(2)} & \dots & a_{2k}^{(2)} & \dots & a_{2n}^{(2)} & x_2 & b_2^{(2)} \\ \vdots & & \vdots & & \vdots & \vdots & \vdots & \vdots \\ a_{kk}^{(k)} & \dots & a_{kn}^{(k)} & \dots & a_{kn}^{(k)} & x_k & b_k^{(k)} & , \\ a_{mk}^{(k)} & \dots & a_{mn}^{(k)} & \dots & a_{mn}^{(k)} & x_n & b_m^{(k)} & \end{array} \quad (2.8)$$

简记为 $\mathbf{A}^{(k)} \mathbf{x} = \mathbf{b}^{(k)}$ 。

设 $a_{kk}^{(k)} \neq 0$, 计算乘数

$$m_{ik} = a_{ik}^{(k)} / a_{kk}^{(k)} \quad (i = k + 1, \dots, m).$$

用 $-m_{ik}$ 乘(2.8)的第 k 个方程加到第 i 个方程($i = k + 1, \dots, m$), 消去从第 $k + 1$ 个方程到第 m 个方程中的未知数 x_k , 得到与(2.1)等价的方程组 $\mathbf{A}^{(k+1)} \mathbf{x} = \mathbf{b}^{(k+1)}$ 。 $\mathbf{A}^{(k+1)}$, $\mathbf{b}^{(k+1)}$ 元素的计算公式为

$$\begin{aligned} a_{ij}^{(k+1)} &= a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)} \quad (i = k + 1, \dots, m; j = k + 1, \dots, n), \\ b_i^{(k+1)} &= b_i^{(k)} - m_{ik} b_k^{(k)} \quad (i = k + 1, \dots, m). \end{aligned} \quad (2.9)$$

显然 $\mathbf{A}^{(k+1)}$ 中从第 1 行到第 k 行与 $\mathbf{A}^{(k)}$ 相同 .

(3) 继续上述过程, 且设 $a_{kk}^{(k)} = 0$ ($i = 1, 2, \dots, s$), 直到完成第 s 步消元计算 . 最后得到与原方程组等价的简单方程组 $\mathbf{A}^{(s+1)} \mathbf{x} = \mathbf{b}^{(s+1)}$, 其中 $\mathbf{A}^{(s+1)}$ 为上梯形 .

特别当 $m = n$ 时, 与原方程组等价的方程组为 $\mathbf{A}^{(n)} \mathbf{x} = \mathbf{b}^{(n)}$, 即

$$\begin{array}{cccccc} a_{11}^{(1)} & a_{21}^{(1)} & \dots & a_{n1}^{(1)} & x_1 & b_1^{(1)} \\ a_{22}^{(2)} & \dots & a_{n2}^{(2)} & x_2 & & b_2^{(2)} \\ \vdots & \dots & \dots & & & \dots \\ a_{nn}^{(n)} & x_n & & & & b_n^{(n)} \end{array} = . \quad (2.10)$$

由(2.1)约化为(2.10)的过程称为消元过程 .

如果 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 是非奇异矩阵, 且 $a_{kk}^{(k)} \neq 0$ ($k = 1, 2, \dots, n - 1$), 求解三角形方程组(2.10), 得到求解公式

$$\begin{aligned} x_n &= b_n^{(n)} / a_{nn}^{(n)}, \\ x_k &= b_k^{(k)} - \sum_{j=k+1}^n a_{kj}^{(k)} x_j / a_{kk}^{(k)} \quad (k = n - 1, n - 2, \dots, 1). \end{aligned} \quad (2.11)$$

(2.10) 的求解过程(2.11)称为回代过程 .

注意: 设 $\mathbf{Ax} = \mathbf{b}$, 其中 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 为非奇异矩阵, 如果 $a_{11} = 0$, 由于 \mathbf{A} 为非奇异矩阵, 所以 \mathbf{A} 的第一列一定有元素不等于零, 例如 $a_{11} \neq 0$, 于是可交换两行元素(即 $r_1 \leftrightarrow r_{i_1}$), 将 a_{11} 调到(1, 1)位置, 然后进行消元计算, 这时 $\mathbf{A}^{(2)}$ 右下角矩阵为 $n - 1$ 阶非奇异矩阵 . 继续这过程, 高斯消去法照样可进行计算 .

总结上述讨论即有

定理 5 设 $\mathbf{Ax} = \mathbf{b}$, 其中 $\mathbf{A} \in \mathbb{R}^{n \times n}$.

(1) 如果 $a_{kk}^{(k)} \neq 0$ ($k = 1, 2, \dots, n$), 则可通过高斯消去法将 $\mathbf{Ax} = \mathbf{b}$ 约化为等价的三角形方程组(2.10), 且计算公式为:

(a) 消元计算($k = 1, 2, \dots, n - 1$)

$$\begin{aligned} m_{ik} &= a_{ik}^{(k)} / a_{kk}^{(k)} \quad (i = k+1, \dots, n), \\ a_{ij}^{(k+1)} &= a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)} \quad (i, j = k+1, \dots, n), \\ b_i^{(k+1)} &= b_i^{(k)} - m_{ik} b_k^{(k)} \quad (i = k+1, \dots, n). \end{aligned}$$

(b) 回代计算

$$\begin{aligned} x_n &= b_n^{(n)} / a_{nn}^{(n)}, \\ x_i &= b_i^{(i)} - \sum_{j=i+1}^n a_{ij}^{(i)} x_j / a_{ii}^{(i)} \quad (i = n-1, \dots, 2, 1). \end{aligned}$$

(2) 如果 \mathbf{A} 为非奇异矩阵, 则可通过高斯消去法(及交换两行的初等变换)将方程组 $\mathbf{Ax} = \mathbf{b}$ 约化为(2.10).

算法1(高斯算法) 设 $\mathbf{A} \in \mathbb{R}^{m \times n}$ ($m > 1$), $s = \min(m-1, n)$, 如果 $a_{kk}^{(k)} \neq 0$ ($k = 1, 2, \dots, s$), 本算法用高斯方法将 \mathbf{A} 约化为上梯形, 且 $\mathbf{A}^{(k)}$ 覆盖 \mathbf{A} , 乘数 m_{ik} 覆盖 a_{ik} .

对于 $k = 1, 2, \dots, s$

(1) 如果 $a_{kk} = 0$, 则计算停止

(2) 对于 $i = k+1, \dots, m$

(a) $a_{ik} \quad m_{ik} = a_{ik} / a_{kk}$

(b) 对于 $j = k+1, \dots, n$

$$a_{ij} \quad a_{ij} - m_{ik} * a_{kj}.$$

显然, 算法1第 k 步需要作 $m-k$ 次除法, $(m-k)(n-k)$ 次乘法运算, 因此, 本算法(从第1步到第 s 步消元计算总的计算量)大约需要 $s^3/3 - (m+n)s^2/2 + mns$ 次乘法运算(对相当大的 s). 当 $m=n$ 时, 总共大约需要 $n^3/3$ 次乘法运算.

数 $a_{kk}^{(k)}$ 在高斯消去法中有着突出的作用, 称为约化的主元素.

算法2(回代算法) 设 $\mathbf{Ux} = \mathbf{b}$, 其中 $\mathbf{U} \in \mathbb{R}^{n \times n}$ 为非奇异上三角阵, 本算法计算 $\mathbf{Ux} = \mathbf{b}$ 的解.

对于 $i = n, \dots, 1$

(1) $x_i = b_i$

(2) 对于 $j = i+1, \dots, n$

$$x_i - \frac{x_i}{a_{ii}} = u_{ij}^* x_j$$

$$(3) \quad x_i - \frac{x_i}{a_{ii}}$$

这个算法需要 $n(n+1)/2$ 乘除法运算.

高斯消去法对于某些简单的矩阵可能会失败, 例如

$$\mathbf{A} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}.$$

由此, 需要对算法 1 进行修改, 首先研究原来矩阵 \mathbf{A} 在什么条件下才能保证 $a_{kk}^{(k)} \neq 0$ ($k=1, 2, \dots$). 下面的定理给出了这个条件.

定理 6 约化的主元素 $a_{ii}^{(i)} \neq 0$ ($i=1, 2, \dots, k$) 的充要条件是矩阵 \mathbf{A} 的顺序主子式 $D_i \neq 0$ ($i=1, 2, \dots, k$) . 即

$$D_1 = a_{11} \neq 0,$$

$$D_i = \begin{vmatrix} a_{11} & \cdots & a_{1i} \\ \cdots & \cdots & \cdots \\ a_{i1} & \cdots & a_{ii} \end{vmatrix} \neq 0 \quad (i=1, 2, \dots, k). \quad (2.12)$$

证明 首先利用归纳法证明定理 6 的充分性. 显然, 当 $k=1$ 时, 定理 6 成立, 现设定理 6 充分性对 $k-1$ 是成立的, 求证定理 6 充分性对 k 亦成立. 设 $D_i \neq 0$ ($i=1, 2, \dots, k$), 于是由归纳法假设 $a_{ii}^{(i)} \neq 0$ ($i=1, 2, \dots, k-1$), 可用高斯消去法将 $\mathbf{A}^{(1)}$ 约化到 $\mathbf{A}^{(k)}$, 即

$$\mathbf{A}^{(1)} = \begin{matrix} a_1^{(1)} & a_2^{(1)} & \cdots & a_{1k}^{(1)} & \cdots & a_n^{(1)} \\ a_2^{(2)} & \cdots & a_{2k}^{(2)} & \cdots & a_n^{(2)} \\ \vdots & & \vdots & & \vdots & \\ a_{nk}^{(k)} & \cdots & a_{nk}^{(k)} & \cdots & a_{nn}^{(k)} & \end{matrix},$$

且有

$$D_2 = \begin{vmatrix} a_{11}^{(1)} & a_{12}^{(1)} \\ 0 & a_{22}^{(2)} \end{vmatrix} = a_{11}^{(1)} a_{22}^{(2)},$$

$$\dots \quad \dots$$

$$D_k = \begin{vmatrix} a_{11}^{(1)} & \dots & a_{1k}^{(1)} \\ \vdots & \ddots & \vdots \\ a_{kk}^{(k)} & \dots & \dots \end{vmatrix} = a_{11}^{(1)} a_{22}^{(2)} \dots a_{kk}^{(k)}. \quad (2.13)$$

由设 $D_i \neq 0$ ($i = 1, 2, \dots, k$), 利用(2.13)式, 则有 $a_{kk}^{(k)} \neq 0$, 定理6充分性对 k 亦成立.

显然, 由假设 $a_{ii}^{(i)} \neq 0$ ($i = 1, 2, \dots, k$), 利用(2.13)式亦可推出 $D_i \neq 0$ ($i = 1, 2, \dots, k$).

推论 如果 \mathbf{A} 的顺序主子式 $D_k \neq 0$ ($k = 1, 2, \dots, n - 1$), 则

$$a_{11}^{(1)} = D_1,$$

$$a_{kk}^{(k)} = D_k / D_{k-1} \quad (k = 2, 3, \dots, n).$$

5.2.2 矩阵的三角分解

下面我们借助矩阵理论进一步对消去法作些分析, 从而建立高斯消去法与矩阵因式分解的关系.

设(2.1)的系数矩阵 $\mathbf{A} \in \mathbf{R}^{n \times n}$ 的各顺序主子式均不为零. 由于对 \mathbf{A} 施行行的初等变换相当于用初等矩阵左乘 \mathbf{A} , 于是对(2.1)施行第一次消元后化为(2.7), 这时 $\mathbf{A}^{(1)}$ 化为 $\mathbf{A}^{(2)}$, $\mathbf{b}^{(1)}$ 化为 $\mathbf{b}^{(2)}$, 即

$$\mathbf{L} \mathbf{A}^{(1)} = \mathbf{A}^{(2)}, \quad \mathbf{L} \mathbf{b}^{(1)} = \mathbf{b}^{(2)},$$

其中

$$\mathbf{L} = \begin{matrix} 1 & & & \\ -m_{21} & 1 & & \\ & \ddots & \ddots & \\ & -m_{n1} & & 1 \end{matrix} \quad \text{W}$$

一般第 k 步消元, $\mathbf{A}^{(k)}$ 化为 $\mathbf{A}^{(k+1)}$, $\mathbf{b}^{(k)}$ 化为 $\mathbf{b}^{(k+1)}$, 相当于

$$\mathbf{L}_k \mathbf{A}^{(k)} = \mathbf{A}^{(k+1)}, \quad \mathbf{L}_k \mathbf{b}^{(k)} = \mathbf{b}^{(k+1)},$$

其中

$$\begin{matrix} & & 1 & & \\ & & w & & \\ \mathbf{L}_k = & & & & \\ & & & 1 & \\ & & & -m_{k+1,k} & 1 \\ & & & \dots & w \\ & & & -m_{nk} & 1 \end{matrix}$$

重复这过程, 最后得到

$$\begin{aligned} \mathbf{L}_{n-1} \dots \mathbf{L}_2 \mathbf{L}_1 \mathbf{A}^{(1)} &= \mathbf{A}^{(n)}; \\ \mathbf{L}_{n-1} \dots \mathbf{L}_2 \mathbf{L}_1 \mathbf{b}^{(1)} &= \mathbf{b}^{(n)}. \end{aligned} \tag{2.14}$$

将上三角矩阵 $\mathbf{A}^{(n)}$ 记为 \mathbf{U} , 由 (2.14) 得到

$$\mathbf{A} = \mathbf{L}^{-1} \mathbf{L}^{-1} \dots \mathbf{L}_{n-1}^{-1} \mathbf{U} = \mathbf{L}\mathbf{U},$$

其中

$$\begin{matrix} & & 1 & & \\ & & m_{11} & 1 & \\ \mathbf{L} = \mathbf{L}^{-1} \mathbf{L}^{-1} \dots \mathbf{L}_{n-1}^{-1} = & & m_{21} & m_{32} & 1 \\ & & \dots & \dots & \dots & w \\ & & m_{n1} & m_{n2} & m_{n3} & \dots & 1 \end{matrix}$$

为单位下三角矩阵.

这就是说, 高斯消去法实质上产生了一个将 \mathbf{A} 分解为两个三角形矩阵相乘的因式分解, 于是我们得到如下重要定理, 它在解方程组的直接法中起着重要作用.

定理 7(矩阵的 LU 分解) 设 \mathbf{A} 为 n 阶矩阵, 如果 \mathbf{A} 的顺序主子式 $D_i \neq 0$ ($i=1, 2, \dots, n-1$), 则 \mathbf{A} 可分解为一个单位下三角矩阵 \mathbf{L} 和一个上三角矩阵 \mathbf{U} 的乘积, 且这种分解是唯一的.

证明 根据以上高斯消去法的矩阵分析, $\mathbf{A} = \mathbf{LU}$ 的存在性已经得到证明, 现仅在 \mathbf{A} 为非奇异矩阵的假定下来证明唯一性, 当 \mathbf{A} 为奇异矩阵的情况留作练习. 设

$$\mathbf{A} = \mathbf{LU} = \mathbf{L} \mathbf{U},$$

其中 \mathbf{L}, \mathbf{L} 为单位下三角矩阵, \mathbf{U}, \mathbf{U} 为上三角矩阵.

由于 \mathbf{U}^{-1} 存在, 故

$$\mathbf{L}^{-1} \mathbf{L} = \mathbf{U} \mathbf{U}^{-1}.$$

上式右边为上三角矩阵, 左边为单位下三角矩阵, 从而上式两边都必须等于单位矩阵, 故 $\mathbf{U} = \mathbf{U}$, $\mathbf{L} = \mathbf{L}$. 证毕.

例 2 对于例 1, 系数矩阵

$$\mathbf{A} = \begin{matrix} 1 & 1 & 1 \\ 0 & 4 & -1 \\ 2 & -2 & 1 \end{matrix},$$

由高斯消去法, $m_{21} = 0$, $m_{31} = 2$, $m_{32} = -1$, 故

$$\mathbf{A} = \begin{matrix} 1 & 0 & 0 & 1 & 1 & 1 \\ 0 & 1 & 0 & 0 & 4 & -1 \\ 2 & -1 & 1 & 0 & 0 & -2 \end{matrix} = \mathbf{LU}.$$

5.3 高斯主元素消去法

由高斯消去法知道, 在消元过程中可能出现 $d_{kk}^{(k)} = 0$ 的情况, 这时消去法将无法进行; 即使主元素 $d_{kk}^{(k)}$ 0 但很小时, 用其作除数, 会导致其他元素数量级的严重增长和舍入误差的扩散, 最后也使得计算解不可靠.

例 3 求解方程组

$$\begin{array}{rrr} 0.001 & 2.000 & 3.000 \\ -1.000 & 3.712 & 4.623 \\ -2.000 & 1.072 & 5.643 \end{array} \begin{array}{l} x_1 \\ x_2 \\ x_3 \end{array} = \begin{array}{r} 1.000 \\ 2.000 \\ 3.000 \end{array}.$$

用 4 位浮点数进行计算 . 精确解舍入到 4 位有效数字为

$$\mathbf{x}^* = (-0.4904, -0.05104, 0.3675)^T,$$

解 (方法 1) 用高斯消去法求解 .

$$(A / b) = \left| \begin{array}{ccc|c} 0.001 & 2.000 & 3.000 & 1.000 \\ -1.000 & 3.712 & 4.623 & 2.000 \\ -2.000 & 1.072 & 5.643 & 3.000 \end{array} \right| \quad \begin{aligned} m_{21} &= -1.000 / 0.001 \\ &= -1000 \\ m_{31} &= -2.000 / 0.001 \\ &= -2000 \end{aligned}$$

$$\left| \begin{array}{ccc|c} 0.001 & 2.000 & 3.000 & 1.000 \\ 0 & 2004 & 3005 & 1002 \\ 0 & 4001 & 6006 & 2003 \end{array} \right| \quad \begin{aligned} m_{32} &= 4001 / 2004 \\ &= 1.997 \end{aligned}$$

$$\left| \begin{array}{ccc|c} 0.001 & 2.000 & 3.000 & 1.000 \\ 0 & 2004 & 3005 & 1002 \\ 0 & 0 & 5.000 & 2.000 \end{array} \right|$$

计算解为

$$\tilde{\mathbf{x}} = (-0.400, -0.09980, 0.4000)^T.$$

显然计算解 $\tilde{\mathbf{x}}$ 是一个很坏的结果, 不能作为方程组的近似解 . 其原因是我们在消元计算时用了小主元 0.001, 使得约化后的方程组元素数量级大大增长, 经再舍入使得在计算(3,3)元素时发生了严重的相消情况 ((3,3)元素舍入到第 4 位数字的正确值是 5.922), 因此经消元后得到的三角形方程组就不准确了 .

(方法 2) 交换行, 避免绝对值小的主元作除数 .

$$(A / b) = \left| \begin{array}{ccc|c} -2.000 & 1.072 & 5.643 & 3.000 \\ -1.000 & 3.712 & 4.623 & 2.000 \\ 0.001 & 2.000 & 3.000 & 1.000 \end{array} \right| \quad \begin{aligned} m_{21} &= 0.5000 \\ m_{31} &= -0.0005 \end{aligned}$$

$$\left| \begin{array}{ccc|c} -2.000 & 1.072 & 5.643 & 3.000 \\ 0 & 3.176 & 1.801 & 0.5000 \\ 0 & 2.001 & 3.003 & 1.002 \end{array} \right| \quad m_{32} = 0.6300$$

$$\left| \begin{array}{ccc|c} -2.000 & 1.072 & 5.643 & 3.000 \\ 0 & 3.176 & 1.801 & 0.5000 \\ 0 & 0 & 1.868 & 0.6870 \end{array} \right|$$

得计算解为

$$\mathbf{x} = (-0.4900, -0.05113, 0.3678)^T \quad \mathbf{x}^*.$$

这个例子告诉我们, 在采用高斯消去法解方程组时, 小主元可能产生麻烦, 故应避免采用绝对值小的主元素 $a_{kk}^{(k)}$. 对一般矩阵来说, 最好每一步选取系数矩阵(或消元后的低阶矩阵)中绝对值最大的元素作为主元素, 以使高斯消去法具有较好的数值稳定性. 这就是全主元素消去法, 在选主元时要花费较多机器时间, 目前主要使用的是列主元消去法, 本节主要介绍列主元消去法, 并假定(2.1)的 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 为非奇异的.

5.3.1 列主元素消去法

设方程组(2.1)的增广矩阵为

$$\mathbf{B} = \left| \begin{array}{cccc|c} a_{11} & a_{12} & \dots & a_{1n} & b_1 \\ a_{21} & a_{22} & \dots & a_{2n} & b_2 \\ \dots & \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} & b_n \end{array} \right|$$

首先在 \mathbf{A} 的第一列中选取绝对值最大的元素作为主元素, 例如

$$|a_{11}| = \max_{1 \leq i \leq n} |a_{ii}| \neq 0,$$

然后交换 \mathbf{B} 的第 1 行与第 i 行, 经第 1 次消元计算得

$$(\mathbf{A} / \mathbf{b}) \rightarrow (\mathbf{A}^{(1)} / \mathbf{b}^{(1)}).$$

重复上述过程, 设已完成第 $k - 1$ 步的选主元素, 交换两行及消元计算, $(\mathbf{A} | \mathbf{b})$ 约化为

$$(A^{(k)} / b^{(k)}) = \begin{array}{cccccc|c} a_{11} & a_{12} & \dots & a_{1k} & \dots & a_{1n} & b_1 \\ a_{21} & \dots & a_{2k} & \dots & a_{2n} & & b_2 \\ \vdots & & \dots & & \dots & & \dots \\ a_{kk} & \dots & a_{kn} & & & & b_k \\ \vdots & & \dots & & & & \dots \\ a_{nk} & \dots & a_{nn} & & & & b_n \end{array},$$

其中 $A^{(k)}$ 的元素仍记为 a_{ij} , $b^{(k)}$ 的元素仍记为 b_i .

第 k 步选主元素(在 $A^{(k)}$ 右下角方阵的第 1 列内选), 即确定 i_k , 使

$$|a_{i_k k}| = \max_{k \leq i \leq n} |a_{ik}| \neq 0.$$

交换($A^{(k)} | b^{(k)}$)第 k 行与 i_k 行的元素, 再进行消元计算, 最后将原方程组化为($k=1, 2, \dots, n-1$)

$$\begin{array}{cccccc|c} a_{11} & a_{12} & \dots & a_{1n} & x_1 & & b \\ a_{21} & \dots & a_{2n} & x_2 & & & b \\ \vdots & & \dots & & & & \dots \\ a_{nn} & x_n & & & & & b_n \end{array}.$$

回代求解

$$x_n = b / a_{nn};$$

$$x_i = b - \sum_{j=i+1}^n a_{ij} x_j / a_{ii} \quad (i = n-1, \dots, 2, 1).$$

算法 3(列主元素消去法) 设 $\mathbf{Ax} = \mathbf{b}$. 本算法用 \mathbf{A} 的具有行交换的列主元素消去法, 消元结果冲掉 \mathbf{A} , 乘数 m_{ij} 冲掉 a_{ij} , 计算解 \mathbf{x} 冲掉常数项 \mathbf{b} , 行列式存放在 \det 中.

1. $\det \rightarrow 1$
2. 对于 $k=1, 2, \dots, n-1$
 - (1) 按列选主元

$$|a_{k,k}| = \max_k |a_{ik}|$$

(2) 如果 $a_{i_k, k} = 0$, 则计算停止 ($\det(\mathbf{A}) = 0$)

(3) 如果 $i_k = k$ 则转(4)

换行: $a_{kj} \quad a_{i_k, j}$ ($j = k, k+1, \dots, n$)

$$b_k \quad b_{i_k}$$

$$\det - \det$$

(4) 消元计算

对于 $i = k+1, \dots, n$

$$(a) a_{ik} \quad m_{ik} = a_{ik} / a_{kk}$$

(b) 对于 $j = k+1, \dots, n$

$$a_{ij} \quad a_{ij} - m_{ik} * a_{kj}$$

$$(c) b_i \quad b_i - m_{ik} * b_k$$

$$(5) \det a_{kk} * \det$$

3. 如果 $a_{nn} = 0$, 则计算停止 ($\det(\mathbf{A}) = 0$)

4. 回代求解

$$(1) b_n \quad b_n / a_{nn}$$

(2) 对于 $i = n-1, \dots, 2, 1$

$$b_i \quad b_i - \sum_{j=i+1}^n a_{ij} * b_j / a_{ii}$$

$$5. \det a_{nn} * \det$$

例 3 的(方法 2)用的就是列主元素消去法.

下面用矩阵运算来描述解(2.1)的列主元素消去法. 列主元素消去法为

$$\begin{aligned} \mathbf{L} \mathbf{I}_{k, i_k} \mathbf{A}^{(1)} &= \mathbf{A}^{(2)}, \quad \mathbf{L} \mathbf{I}_{k, i_k} \mathbf{b}^{(1)} = \mathbf{b}^{(2)}, \\ \mathbf{L} \mathbf{I}_{k, i_k} \mathbf{A}^{(k)} &= \mathbf{A}^{(k+1)}, \quad \mathbf{L} \mathbf{I}_{k, i_k} \mathbf{b}^{(k)} = \mathbf{b}^{(k+1)}. \end{aligned} \quad (3.1)$$

其中 \mathbf{L}_k 的元素满足 $|m_{ik}| \leq 1$ ($k = 1, 2, \dots, n-1$), \mathbf{I}_{k, i_k} 是初等置换阵.

利用(3.1)得到

$$\mathbf{L}_{n-1} \mathbf{L}_{n-1, i_{n-1}} \dots \mathbf{L}_2 \mathbf{L}_{i_2} \mathbf{L}_{i_1} \mathbf{A} = \mathbf{A}^{(n)} = \mathbf{U}.$$

简记为

$$\mathbf{R}\mathbf{A} = \mathbf{U}, \mathbf{R}\mathbf{b} = \mathbf{b}^{(n)},$$

其中

$$\mathbf{R} = \mathbf{L}_{n-1} \mathbf{L}_{n-1, i_{n-1}} \dots \mathbf{L}_2 \mathbf{L}_{i_2} \mathbf{L}_{i_1}.$$

下面就 $n=4$ 来考察一下矩阵 \mathbf{R} .

$$\begin{aligned} \mathbf{U} = \mathbf{A}^{(4)} &= \mathbf{L}_{i_3} \mathbf{L}_{i_3, i_3} \mathbf{L}_{i_2} \mathbf{L}_{i_2, i_2} \mathbf{L}_{i_1} \mathbf{A} \\ &= \mathbf{L}_{i_3} (\mathbf{I}_{i_3} \mathbf{L}_{i_2} \mathbf{I}_{i_3}) (\mathbf{I}_{i_3} \mathbf{L}_{i_2, i_2} \mathbf{L}_{i_2} \mathbf{I}_{i_3}) \\ &\quad \cdot (\mathbf{I}_{i_3} \mathbf{L}_{i_2, i_2} \mathbf{I}_{i_1}) \mathbf{A} \quad \text{即 } \mathbf{R} = \mathbf{PA}, \end{aligned} \quad (3.2)$$

其中

$$\mathbf{R}_1 = \mathbf{I}_{i_3} \mathbf{L}_{i_2} \mathbf{L}_{i_2, i_2} \mathbf{I}_{i_3},$$

$$\mathbf{R}_2 = \mathbf{I}_{i_3} \mathbf{L}_{i_2} \mathbf{I}_{i_3},$$

$$\mathbf{R}_3 = \mathbf{L},$$

$$\mathbf{P} = \mathbf{I}_{i_3} \mathbf{L}_{i_2} \mathbf{I}_{i_1}.$$

由习题 3 知 \mathbf{R}_k ($k=1, 2, 3$) 亦为单位下三角阵, 其元素的绝对值不超过 1. 记

$$\mathbf{L}^{-1} = \mathbf{R}_3 \mathbf{R}_2 \mathbf{R}_1,$$

由(3.2)得到

$$\mathbf{PA} = \mathbf{LU},$$

其中 \mathbf{P} 为排列矩阵, \mathbf{L} 为单位下三角阵, \mathbf{U} 为上三角阵. 这说明对(2.1)应用列主元素消去法相当于对 $(\mathbf{A} | \mathbf{b})$ 先进行一系列行交换后对 $\mathbf{PAx} = \mathbf{Pb}$ 再应用高斯消去法. 在实际计算中我们只能在计算过程中做行的交换.

总结以上的讨论有

定理 8(列主元素的三角分解定理) 如果 \mathbf{A} 为非奇异矩阵,

则存在排列矩阵 \mathbf{P} 使

$$\mathbf{PA} = \mathbf{LU},$$

其中 \mathbf{L} 为单位下三角阵, \mathbf{U} 为上三角阵 .

在编程实现过程中, \mathbf{L} 元素存放在数组 \mathbf{A} 的下三角部分, \mathbf{U} 元素存放在 \mathbf{A} 上三角部分, 由记录主行的整型数组 $\text{Ip}(n)$ 可知 \mathbf{P} 的情况 .

5.3.2 高斯 - 若当消去法

高斯消去法始终是消去对角线下方的元素, 现考虑高斯消去法的一种修正, 即消去对角线下方和上方的元素, 这种方法称为高斯 - 若当 (Gauss-Jordan) 消去法 .

设用高斯 - 若当消去法已完成 $k - 1$ 步, 于是 $\mathbf{Ax} = \mathbf{b}$ 化为等价方程组 $\mathbf{A}^{(k)} \mathbf{x} = \mathbf{b}^{(k)}$, 其中

$$(\mathbf{A}^{(k)} / \mathbf{b}^{(k)}) = \left| \begin{array}{cccccc|c} 1 & 0 & \dots & 0 & a_{1k} & \dots & a_{1n} & b_1 \\ 1 & \dots & 0 & a_{kk} & \dots & a_{kn} & b_2 \\ \vdots & & \dots & & & \dots & \dots \\ 1 & & & a_{k-1,k} & \dots & a_{k-1,n} & b_{k-1} \\ a_{kk} & & & \dots & & a_{kn} & b_k \\ \dots & & & & & \dots & \dots \\ a_{nk} & & & a_{nn} & & & b_n \end{array} \right|.$$

在第 k 步计算时 ($k = 1, 2, \dots, n$), 考虑对上述矩阵的第 k 行上、下都进行消元计算 .

1. 按列选主元素, 即确定 i_k 使

$$|a_{i_k k}| = \max_{k \leq i \leq n} |a_{ik}|.$$

2. 换行 (当 $i_k \neq k$ 时) 交换 ($\mathbf{A} | \mathbf{b}$) 第 k 行与第 i_k 行元素 .

3. 计算乘数 $m_{ik} = -a_{ik}/a_{kk}$ ($i = 1, 2, \dots, n$ 且 $i \neq k$),

$$m_{kk} = 1/a_{kk}.$$

(m_{ik} 可保存在存放 a_{ik} 的单元中) .

4. 消元计算

$$\begin{aligned} a_{ij} &= a_{ij} + m_{ik} a_{kj} \quad i = 1, 2, \dots, n \text{ 且 } i \neq k \\ b_i &= b_i + m_{ik} b_k \quad (i = 1, 2, \dots, n \text{ 且 } i \neq k) . \end{aligned}$$

5. 计算主行

$$\begin{aligned} a_{kj} &= a_{kj} / m_{kk} \quad (j = k, k+1, \dots, n) , \\ b_k &= b_k / m_{kk} . \end{aligned}$$

上述过程结束后有

$$(A / \mathbf{b}) \rightarrow (A^{(k+1)} / \mathbf{b}^{(k+1)}) = \begin{array}{c|c} 1 & \mathbf{b}_1 \\ 1 & \mathbf{b}_2 \\ \vdots & \vdots \\ 1 & \mathbf{b}_n \end{array} .$$

说明用高斯-若当方法将 A 约化为单位矩阵, 计算解就在常数项位置得到, 因此用不着回代求解, 用高斯-若当方法解方程组其计算量大约需要 $n^3/2$ 次乘除法, 要比高斯消去法大, 但用高斯-若当方法求一个矩阵的逆矩阵还是比较合适的.

定理 9(高斯-若当法求逆矩阵) 设 A 为非奇异矩阵, 方程组 $AX = I_n$ 的增广矩阵为 $C = (A | I_n)$. 如果对 C 应用高斯-若当方法化为 $(I_n | T)$, 则 $A^{-1} = T$.

事实上, 求 A 的逆矩阵 A^{-1} , 即求 n 阶矩阵 X , 使 $AX = I_n$, 其中 I_n 为单位矩阵. 将 X 按列分块

$$X = (\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_n), \quad I = (\mathbf{e}_1 \ \mathbf{e}_2 \ \dots \ \mathbf{e}_n),$$

于是求解 $AX = I$ 等价于求解 n 个方程组

$$A\mathbf{x}_j = \mathbf{e}_j \quad (j = 1, 2, \dots, n) .$$

我们可用高斯-若当方法求解 $AX = I$.

例 4 用高斯-若当方法求

$$\mathbf{A} = \begin{matrix} & 1 & 2 & 3 \\ & 2 & 4 & 5 \\ & 3 & 5 & 6 \end{matrix}$$

的逆矩阵 \mathbf{A}^{-1} .

$$\text{解 } \mathbf{C} = \left| \begin{array}{ccc|ccc|ccc} 1 & 2 & 3 & 1 & 0 & 0 & 3 & 5 & 6 & 0 & 0 & 1 \\ 2 & 4 & 5 & 0 & 1 & 0 & r_1 & & 2 & 4 & 5 & 0 & 1 & 0 \\ 3 & 5 & 6 & 0 & 0 & 1 & & & 1 & 2 & 3 & 1 & 0 & 0 \end{array} \right.$$

$$\text{第 1 次消元} \quad \left| \begin{array}{ccc|ccc} 1 & 5/3 & 2 & 0 & 0 & 1/3 \\ 0 & 2/3 & 1 & 0 & 1 & -2/3 \\ 0 & 1/3 & 1 & 1 & 0 & -1/3 \end{array} \right.$$

c

$$\text{第 2 次消元} \quad \left| \begin{array}{ccc|ccc} 1 & 0 & -1/2 & 0 & -5/2 & 2 \\ 0 & 1 & 3/2 & 0 & 3/2 & -1 \\ 0 & 0 & 1/2 & 1 & -1/2 & 0 \end{array} \right.$$

c

$$\text{第 3 次消元} \quad \left| \begin{array}{ccc|ccc} 1 & 0 & 0 & 1 & -3 & 2 \\ 0 & 1 & 0 & -3 & 3 & -1 \\ 0 & 0 & 1 & 2 & -1 & 0 \end{array} \right. = (\mathbf{I}_3 | \mathbf{A}^{-1}) .$$

c

且 $\mathbf{m} = (m_{11}, m_{21}, m_{31})^T = \mathbf{c}$, $\mathbf{m} = (m_{12}, m_{22}, m_{32})^T = \mathbf{c}$,
 $\mathbf{m} = (m_{13}, m_{23}, m_{33})^T = \mathbf{c}$.

为了节省内存单元, 可不必将单位矩阵存放起来, **c** 存放在 **A** 的第 1 列位置, **c** 存放在 **A** 的第二列位置, **c** 存放在 **A** 的第 3 列位置, 经消元计算, 最后再调整一下列就可在 **A** 的位置得到 **A**⁻¹. 注意第 k 步消元时, 由 **A** 的第 k 列

$$\mathbf{a} = (a_{1k}, \dots, a_{kk}, \dots, a_{nk})^T$$

计算

$$\mathbf{m} = -\frac{a_{kk}}{a_{kk}}, \dots, \frac{1}{a_{kk}}, \dots, -\frac{a_{kk}}{a_{kk}}^T,$$

且冲掉 \mathbf{a} .

最后在 \mathbf{A} 位置如何调整列呢! 事实上, 我们在 \mathbf{A} 位置最后得到矩阵 $\mathbf{PA} = \mathbf{A}$ (其中 \mathbf{P} 为排列阵) 的逆矩阵 \mathbf{A}^{-1} , 于是

$$\mathbf{A}^{-1} = \mathbf{A}^{-1} \mathbf{P}.$$

5.4 矩阵三角分解法

高斯消去法有很多变形, 有的是高斯消去法的改进、改写, 有的是用于某一类特殊性质矩阵的高斯消去法的简化.

5.4.1 直接三角分解法

将高斯消去法改写为紧凑形式, 可以直接从矩阵 \mathbf{A} 的元素得到计算 \mathbf{L}, \mathbf{U} 元素的递推公式, 而不需任何中间步骤, 这就是所谓直接三角分解法. 一旦实现了矩阵 \mathbf{A} 的 LU 分解, 那么求解 $\mathbf{Ax} = \mathbf{b}$ 的问题就等价于求解两个三角形方程组

$$\mathbf{Ly} = \mathbf{b}, \text{ 求 } \mathbf{y};$$

$$\mathbf{Ux} = \mathbf{y}, \text{ 求 } \mathbf{x}.$$

1. 不选主元的三角分解法

设 \mathbf{A} 为非奇异矩阵, 且有分解式

$$\mathbf{A} = \mathbf{LU},$$

其中 \mathbf{L} 为单位下三角阵, \mathbf{U} 为上三角阵, 即

$$\mathbf{A} = \begin{matrix} & 1 & & & & \\ & l_{21} & 1 & & & \\ & \cdots & \cdots & \mathbf{W} & & \\ l_{n1} & l_{n2} & \cdots & 1 & & \end{matrix} \begin{matrix} u_{11} & u_{12} & \cdots & u_{1n} \\ u_{22} & \cdots & u_{2n} \\ \mathbf{W} & \cdots & \\ & & u_{nn} \end{matrix}. \quad (4.1)$$

下面说明 \mathbf{L}, \mathbf{U} 的元素可以由 n 步直接计算定出, 其中第 r 步定出 \mathbf{U} 的第 r 行和 \mathbf{L} 的第 r 列元素. 由(4.1)有:

$$a_{1i} = u_{1i} \quad (i = 1, 2, \dots, n), \text{ 得 } \mathbf{U} \text{ 的第 1 行元素;}$$

$$a_{ii} = l_{11} u_{11}, \quad l_{11} = a_{ii}/u_{11} \quad (i = 2, \dots, n), \text{ 得 } \mathbf{L} \text{ 的第 1 列元素.}$$

设已经定出 \mathbf{U} 的第 1 行到第 $r-1$ 行元素与 \mathbf{L} 的第 1 列到第 $r-1$ 列元素. 由(4.1), 利用矩阵乘法(注意当 $r < k$ 时, $l_{rk} = 0$), 有

$$a_{ri} = \sum_{k=1}^n l_{rk} u_{ki} = \sum_{k=1}^{r-1} l_{rk} u_{ki} + u_{ri}.$$

故

$$u_{ri} = a_{ri} - \sum_{k=1}^{r-1} l_{rk} u_{ki} \quad (i = r, r+1, \dots, n),$$

又由(4.1)有

$$a_{ir} = \sum_{k=1}^n l_{ik} u_{kr} = \sum_{k=1}^{r-1} l_{ik} u_{kr} + l_{ir} u_{rr}.$$

总结上述讨论, 得到用直接三角分解法解 $\mathbf{Ax} = \mathbf{b}$ (要求 \mathbf{A} 的所有顺序主子式都不为零)的计算公式.

$$u_{1i} = a_{1i} \quad (i = 1, 2, \dots, n), \quad l_{11} = a_{11}/u_{11} \quad (i = 2, 3, \dots, n),$$

计算 \mathbf{U} 的第 r 行, \mathbf{L} 的第 r 列元素($r = 2, 3, \dots, n$).

$$u_{ri} = a_{ri} - \sum_{k=1}^{r-1} l_{rk} u_{ki} \quad (i = r, r+1, \dots, n); \quad (4.2)$$

$$l_{ir} = a_{ir} - \sum_{k=1}^{r-1} l_{ik} u_{kr} / u_{rr} \quad (i = r+1, \dots, n; \text{ 且 } r < n); \quad (4.3)$$

求解 $\mathbf{Ly} = \mathbf{b}$, $\mathbf{Ux} = \mathbf{y}$ 的计算公式;

$$\begin{aligned} y_1 &= b; \\ y_i &= b - \sum_{k=1}^{i-1} l_{ik} y_k \quad (i = 2, 3, \dots, n); \end{aligned} \quad (4.4)$$

$$\begin{aligned}
 x_n &= y_n / u_{nn}; \\
 x_i &= y_i - \sum_{k=i+1}^n u_{ik} x_k / u_{ii} \quad (i = n-1, n-2, \dots, 1).
 \end{aligned} \tag{4.5}$$

例 5 用直接三角分解法解

$$\begin{array}{ccccc}
 1 & 2 & 3 & x_1 & 14 \\
 2 & 5 & 2 & x_2 & = 18 \\
 3 & 1 & 5 & x_3 & 20
 \end{array}.$$

解 用分解公式(4.2)—(4.3)计算得

$$\mathbf{A} = \begin{array}{ccccc}
 1 & 0 & 0 & 1 & 2 & 3 \\
 2 & 1 & 0 & 0 & 1 & -4 \\
 3 & -5 & 1 & 0 & 0 & -24
 \end{array} = \mathbf{LU}.$$

求解

$$\begin{aligned}
 \mathbf{Ly} &= (14, 18, 20)^T, \text{ 得 } \mathbf{y} = (14, -10, -72)^T, \\
 \mathbf{Ux} &= (14, -10, -72)^T, \text{ 得 } \mathbf{x} = (1, 2, 3)^T.
 \end{aligned}$$

由于在计算机实现时当 u_{ri} 计算好后 a_{ri} 就不用了, 因此计算好 \mathbf{L}, \mathbf{U} 的元素后就存放在 \mathbf{A} 的相应位置. 例如

$$\mathbf{A} = \begin{array}{cccc|cccc}
 a_{11} & a_{12} & a_{13} & a_{14} & u_{11} & u_{12} & u_{13} & u_{14} \\
 a_{21} & a_{22} & a_{23} & a_{24} & l_{11} & u_{22} & u_{23} & u_{24} \\
 a_{31} & a_{32} & a_{33} & a_{34} & l_{12} & l_{23} & u_{33} & u_{34} \\
 a_{41} & a_{42} & a_{43} & a_{44} & l_{13} & l_{24} & l_{34} & u_{44}
 \end{array}$$

最后在存放 \mathbf{A} 的数组中得到 \mathbf{L}, \mathbf{U} 的元素.

由直接三角分解计算公式, 需要计算形如 $a_i b_i$ 的式子, 可采用“双精度累加”, 以提高精度.

直接分解法大约需要 $n^3/3$ 次乘除法, 和高斯消去法计算量基本相同.

如果已经实现了 $\mathbf{A} = \mathbf{LU}$ 的分解计算, 且 \mathbf{L}, \mathbf{U} 保存在 \mathbf{A} 的相应位置, 则用直接三角分解法解具有相同系数的方程组 $\mathbf{Ax} = (\mathbf{b}_1 \mathbf{b}_2 \dots \mathbf{b}_n)$ 是相当方便的, 每解一个方程组 $\mathbf{Ax} = \mathbf{b}_i$ 仅需要增加 n^2 次乘除法运算.

矩阵 \mathbf{A} 的分解公式(4.2), (4.3)又称为杜利特尔(Doolittle)分解.

2. 选主元的三角分解法

从直接三角分解公式可看出当 $u_{rr} = 0$ 时计算将中断, 或者当 u_{rr} 绝对值很小时, 按分解公式计算可能引起舍入误差的累积. 但如果 \mathbf{A} 非奇异, 我们可通过交换 \mathbf{A} 的行实现矩阵 \mathbf{PA} 的 LU 分解, 因此可采用与列主元消去法类似的方法(可以证明下述方法与列主元消去法等价), 将直接三角分解法修改为(部分)选主元的三角分解法.

设第 $r - 1$ 步分解已完成, 这时有

$$\begin{array}{ccccccc} u_{11} & u_{12} & \dots & u_{1, r-1} & u_{1r} & \dots & u_{1n} \\ b_1 & b_2 & \dots & b_{r-1} & b_r & \dots & b_n \\ \dots & \dots & & \dots & \dots & & \dots \\ \mathbf{A} & l_{r-1,1} & l_{r-1,2} & \dots & u_{r-1,r-1} & u_{r-1,r} & \dots & u_{r-1,n} \\ l_{r1} & l_{r2} & \dots & l_{r,r-1} & a_{rr} & \dots & a_{rn} \\ \dots & \dots & & \dots & \dots & & \dots \\ l_{n1} & l_{n2} & \dots & l_{n,r-1} & a_{nr} & \dots & a_{nn} \end{array}$$

第 r 步分解需用到(4.2)及(4.3)式, 为了避免用小的数 u_{rr} 作除数, 引进量

$$s_i = a_{ir} - \sum_{k=1}^{r-1} l_{ik} u_{kr} \quad (i = r, r+1, \dots, n).$$

于是有

$$u_{rr} = s_r, \quad l_{ir} = s_r / s_r \quad (i = r+1, \dots, n).$$

取 $\max_{r \leq i \leq n} |s_i| = |s_r|$, 交换 \mathbf{A} 的 r 行与 i_r 行元素, 将 s_{i_r} 调到 (r, r) 位置(将 (i, j) 位置的新元素仍记为 l_{ij} 及 a_{ij}), 于是有 $|l_{ir}| = 1$ ($i = r+1, \dots, n$) . 由此再进行第 r 步分解计算 .

算法 4(选主元的三角分解法) 设 $\mathbf{Ax} = \mathbf{b}$, 其中 \mathbf{A} 为非奇异矩阵 . 本算法采用选主元的三角分解法, 用 $\mathbf{PA} = \mathbf{L}_{n-1, i_{n-1}} \dots \mathbf{L}_{1, i_1} \mathbf{A}$ 的三角分解冲掉 \mathbf{A} , 用整型数组 $\text{Ip}(n)$ 记录主行, 解 \mathbf{x} 存放在 \mathbf{b} 内 .

1. 对于 $r=1, 2, \dots, n$

(1) 计算 s_i

$$a_{ir} - s_i = a_{ir} - \sum_{k=1}^{r-1} l_{ik} u_{kr} \quad (i = r, r+1, \dots, n)$$

(2) 选主元 $|s_r| = \max_{r \leq i \leq n} |s_i|$, $\text{Ip}(r) = i_r$

(3) 交换 \mathbf{A} 的 r 行与 i_r 行元素

$$a_{ri} = a_{i_r i} \quad (i = 1, 2, \dots, n)$$

(4) 计算 \mathbf{U} 的第 r 行元素, \mathbf{L} 的第 r 列元素

$$a_{rr} = u_{rr} = s_r$$

$$a_{ir} - l_{ir} = s_r / u_{rr} = a_{ir} / a_{rr} \quad (i = r+1, \dots, n, \text{且 } r \neq n)$$

$$a_{ri} - u_{ri} = a_{ri} - \sum_{k=1}^{r-1} l_{rk} u_{ki} \quad (i = r+1, \dots, n, \text{且 } r \neq n)$$

(这时有 $|l_{ir}| = 1$) .

上述计算过程完成后就实现了 \mathbf{PA} 的 LU 分解, 且 \mathbf{U} 保存在 \mathbf{A} 的上三角部分, \mathbf{L} 保存在 \mathbf{A} 的下三角部分, 排列阵 \mathbf{P} 由 $\text{Ip}(n)$ 最后记录可知 .

求解 $\mathbf{Ly} = \mathbf{Pb}$ 及 $\mathbf{Ux} = \mathbf{y}$.

2. 对于 $i=1, \dots, n-1$

(1) $t = \text{Ip}(i)$

(2) 如果 $i=t$ 则转(3)

$$b_t = b_i$$

(3) (继续循环)

$$3. b_i - \sum_{k=1}^{i-1} l_{ik} b_k \quad (i = 2, 3, \dots, n)$$

$$4. b_n - b / u_{nn}, b_i - \sum_{k=i+1}^n u_{ik} b_k / u_{ii} \quad (i = n-1, \dots, 1)$$

利用算法 4 的结果(实现 $\mathbf{PA} = \mathbf{LU}$ 三角分解), 则可以计算 \mathbf{A} 的逆矩阵

$$\mathbf{A}^{-1} = \mathbf{U}^{-1} \mathbf{L}^{-1} \mathbf{P}.$$

利用 \mathbf{PA} 的三角分解计算 \mathbf{A}^{-1} 步骤:

- (1) 计算上三角矩阵的逆阵 \mathbf{U}^{-1} ;
- (2) 计算 $\mathbf{U}^{-1} \mathbf{L}^{-1}$;
- (3) 交换 $\mathbf{U}^{-1} \mathbf{L}^{-1}$ 列(利用 $\text{Ip}(n)$ 最后记录).

上述方法求 \mathbf{A}^{-1} 大约需要 n^3 次乘法运算.

5.4.2 平方根法

应用有限元法解结构力学问题时, 最后归结为求解线性方程组, 系数矩阵大多具有对称正定性质. 所谓平方根法, 就是利用对称正定矩阵的三角分解而得到的求解对称正定方程组的一种有效方法, 目前在计算机上广泛应用平方根法解此类方程组.

设 \mathbf{A} 为对称矩阵, 且 \mathbf{A} 的所有顺序主子式均不为零, 由本章定理 7 知, \mathbf{A} 可唯一分解为如(4.1)的形式.

为了利用 \mathbf{A} 的对称性, 将 \mathbf{U} 再分解为

$$\mathbf{U} = \begin{matrix} & \frac{u_{12}}{u_{11}} & \cdots & \cdots & \frac{u_{1n}}{u_{11}} \\ u_{11} & 1 & & & \\ & \frac{u_{23}}{u_{22}} & \cdots & \cdots & \frac{u_{2n}}{u_{22}} \\ & u_{22} & 1 & & \\ & & \frac{u_{34}}{u_{33}} & \cdots & \cdots \\ & & u_{33} & 1 & \\ & & & \ddots & \end{matrix} = \mathbf{DU}_0,$$

其中 \mathbf{D} 为对角阵, \mathbf{U} 为单位上三角阵. 于是

$$\mathbf{A} = \mathbf{L}\mathbf{U} = \mathbf{LDU}_0. \quad (4.6)$$

又

$$\mathbf{A} = \mathbf{A}^T = \mathbf{U}_0^T (\mathbf{DL}^T),$$

由分解的唯一性即得

$$\mathbf{U}_0^T = \mathbf{L}.$$

代入(4.6)得到对称矩阵 \mathbf{A} 的分解式 $\mathbf{A} = \mathbf{LDL}^T$. 总结上述讨论有

定理 10(对称阵的三角分解定理) 设 \mathbf{A} 为 n 阶对称阵, 且 \mathbf{A} 的所有顺序主子式均不为零, 则 \mathbf{A} 可唯一分解为

$$\mathbf{A} = \mathbf{LDL}^T,$$

其中 \mathbf{L} 为单位下三角阵, \mathbf{D} 为对角阵.

现设 \mathbf{A} 为对称正定矩阵. 首先说明 \mathbf{A} 的分解式 $\mathbf{A} = \mathbf{LDL}^T$ 中 \mathbf{D} 的对角元素 d_i 均为正数.

事实上, 由 \mathbf{A} 的对称正定性, 5.2 节中的推论成立, 即

$$d_1 = D_1 > 0, \quad d_i = D_i / D_{i-1} > 0 \quad (i = 2, 3, \dots, n).$$

于是

$$\begin{aligned} \mathbf{D} &= \begin{matrix} d_1 & & & \\ & \ddots & & \\ & & d_i & \\ & & & \ddots & d_n \end{matrix} \\ &= \mathbf{D}^{\frac{1}{2}} \mathbf{D}^{\frac{1}{2}}, \end{aligned}$$

由定理 6 得到

$$\begin{aligned} \mathbf{A} = \mathbf{LDL}^T &= \mathbf{LD}^{\frac{1}{2}} \mathbf{D}^{\frac{1}{2}} \mathbf{L}^T = (\mathbf{LD}^{\frac{1}{2}})(\mathbf{LD}^{\frac{1}{2}})^T \\ &= \mathbf{L} \mathbf{L}^T, \end{aligned}$$

其中 $\mathbf{L} = \mathbf{LD}^{\frac{1}{2}}$ 为下三角矩阵.

定理 11(对称正定矩阵的三角分解或 Cholesky 分解) 如果 \mathbf{A} 为 n 阶对称正定矩阵, 则存在一个实的非奇异下三角阵 \mathbf{L} 使

$\mathbf{A} = \mathbf{L}\mathbf{L}^T$, 当限定 \mathbf{L} 的对角元素为正时, 这种分解是唯一的.

下面我们用直接分解方法来确定计算 \mathbf{L} 元素的递推公式.
因为

$$\mathbf{A} = \begin{array}{ccccccccc} l_{11} & & & l_{11} & l_{21} & \dots & l_{n1} \\ b_1 & b_2 & & l_{22} & \dots & l_{n2} \\ \dots & \dots & w & w & \dots & \\ l_{n1} & l_{n2} & \dots & l_{nn} & & l_{nn} \end{array},$$

其中 $l_{ii} > 0 (i = 1, 2, \dots, n)$. 由矩阵乘法及 $l_{jk} = 0$ (当 $j < k$ 时), 得

$$a_{ij} = \sum_{k=1}^n l_{ik} l_{jk} = \sum_{k=1}^{j-1} l_{ik} l_{jk} + l_{jj} l_{ij},$$

于是得到解对称正定方程组 $\mathbf{Ax} = \mathbf{b}$ 的平方根法计算公式:

对于 $j = 1, 2, \dots, n$

$$\begin{aligned} 1. \quad l_{ij} &= a_{ij} - \sum_{k=1}^{j-1} \hat{l}_{jk}^2, \\ 2. \quad l_{ij} &= a_{ij} - \sum_{k=1}^{j-1} l_{ik} l_{jk} / l_{jj} \quad (i = j+1, \dots, n); \end{aligned} \tag{4.7}$$

求解 $\mathbf{Ax} = \mathbf{b}$, 即求解两个三角形方程组:

(1) $\mathbf{Ly} = \mathbf{b}$, 求 \mathbf{y} ; (2) $\mathbf{L}^T \mathbf{x} = \mathbf{y}$, 求 \mathbf{x} .

$$\begin{aligned} 3. \quad y_i &= b_i - \sum_{k=1}^n l_{ki} y_k / l_{ii} \quad (i = 1, 2, \dots, n). \\ 4. \quad x_i &= b_i - \sum_{k=i+1}^n l_{ki} x_k / l_{ii} \quad (i = n, n-1, \dots, 1). \end{aligned} \tag{4.8}$$

由计算公式 1 知

$$a_{jj} = \hat{l}_{jj} \quad (j = 1, 2, \dots, n),$$

所以

$$\hat{l}_{jk} = a_{jj} \max_{1 \leq j \leq n} \{a_{jj}\},$$

于是

$$\max_{j,k} \{ l_{jk}^2 \} = \max_{1 \leq j \leq n} \{ a_{jj} \} .$$

上面分析说明, 分解过程中元素 l_{jk} 的数量级不会增长且对角元素 l_{jj} 恒为正数. 于是不选主元素的平方根法是一个数值稳定的方法.

当求出 \mathbf{L} 的第 j 列元素时, \mathbf{L}^T 的第 j 行元素亦算出. 所以平方根法约需 $n^3/6$ 次乘除法, 大约为一般直接 \mathbf{LU} 分解法计算量的一半.

由于 \mathbf{A} 为对称阵, 因此在计算机实现时只需存储 \mathbf{A} 的下三角部分, 共需要存储 $n(n+1)/2$ 个元素, 可用一维数组存放, 即

$$A(n \cdot (n+1)/2) = \{ a_{11}, a_{21}, a_{22}, \dots, a_{n1}, a_{n2}, \dots, a_{nn} \} .$$

矩阵元素 a_{ij} 一维数组的表示为 $A(i \cdot (i-1)/2 + j)$, \mathbf{L} 的元素存放在 A 的相应位置.

由公式(4.7)看出, 用平方根法解对称正定方程组时, 计算 \mathbf{L} 的元素 l_{ii} 需要用到开方运算. 为了避免开方, 我们下面用定理 10 的分解式

$$\mathbf{A} = \mathbf{LDL}^T ,$$

即

$$\mathbf{A} = \begin{array}{ccccccccc} 1 & & d_1 & & 1 & l_{11} & \dots & l_{n1} \\ l_{21} & 1 & & d_2 & & 1 & \dots & l_{n2} \\ \dots & \dots & w & & w & & w & \dots \\ l_{n1} & l_{n2} & \dots & 1 & & d_n & & 1 \end{array} .$$

由矩阵乘法, 并注意 $l_{jj} = 1$, $l_{jk} = 0(j < k)$, 得

$$\begin{aligned} a_{ij} &= \sum_{k=1}^n (\mathbf{LD})_{ik} (\mathbf{L}^T)_{kj} = \sum_{k=1}^n l_{ik} d_k l_{kj} \\ &= \sum_{k=1}^{j-1} l_{ik} d_k l_{kj} + l_{ij} d_j l_{jj} . \end{aligned}$$

于是得到计算 \mathbf{L} 的元素及 \mathbf{D} 的对角元素公式:

对于 $i = 1, 2, \dots, n$.

$$\begin{aligned} 1. \quad l_{ij} &= a_{ij} - \sum_{k=1}^{j-1} l_{ik} d_k l_{jk} / d_j \quad (j = 1, 2, \dots, i-1); \quad (4.9) \\ 2. \quad d_i &= a_{ii} - \sum_{k=1}^{i-1} l_{ik}^2 d_k. \end{aligned}$$

为了避免重复计算, 我们引进

$$t_{ij} = l_{ij} d_j,$$

由(4.9)得到按行计算 \mathbf{L}, \mathbf{T} 元素的公式:

$$d_1 = a_{11}$$

对于 $i = 2, 3, \dots, n$.

$$\begin{aligned} 1. \quad t_{ij} &= a_{ij} - \sum_{k=1}^{j-1} t_{ik} l_{jk}, \quad (j = 1, 2, \dots, i-1); \\ 2. \quad l_{ij} &= t_{ij} / d_j \quad (j = 1, 2, \dots, i-1); \\ 3. \quad d_i &= a_{ii} - \sum_{k=1}^{i-1} t_{ik} l_{ik}. \end{aligned} \quad (4.10)$$

计算出 $\mathbf{T} = \mathbf{LD}$ 的第 i 行元素 t_{ij} ($j = 1, 2, \dots, i-1$) 后, 存放在 \mathbf{A} 的第 i 行相应位置, 然后再计算 \mathbf{L} 的第 i 行元素, 存放在 \mathbf{A} 的第 i 行. \mathbf{D} 的对角元素存放在 \mathbf{A} 的相应位置. 例如

$$\mathbf{A} = \begin{array}{ccccc} a_{11} & & & & d_1 \\ a_{21} & a_{22} & \text{对称} & & d_2 \\ a_{31} & a_{32} & a_{33} & & d_3 \\ a_{41} & a_{42} & a_{43} & a_{44} & d_4 \\ \hline d_1 & & & & \\ l_1 & d_2 & & & \\ l_1 & l_2 & d_3 & & \\ l_1 & l_2 & l_3 & d_4 & \end{array}.$$

对称正定矩阵 \mathbf{A} 按 \mathbf{LDL}^T 分解和按 \mathbf{LL}^T 分解计算量差不多,

但 LDL^T 分解不需要开方计算 .

求解 $\mathbf{Ly} = \mathbf{b}$, $\mathbf{DL}^T \mathbf{x} = \mathbf{y}$ 计算公式

$$4. \quad y_1 = b; \quad (4.11)$$

$$y_i = b - \sum_{k=1}^{i-1} l_{ki} y_k \quad (i = 2, \dots, n).$$

$$5. \quad x_n = y_n / d_n;$$

$$x_i = y_i / d_i - \sum_{k=i+1}^n l_{ki} x_k \quad (i = n-1, \dots, 2, 1).$$

计算公式(4.10), (4.11)称为改进的平方根法 .

5.4.3 追赶法

在一些实际问题中, 例如解常微分方程边值问题, 解热传导方程以及船体数学放样中建立三次样条函数等, 都会要求解系数矩阵为对角占优的三对角线方程组

$$\begin{array}{ccccccccc} b & c_1 & & & x_1 & & f_1 \\ a & b & c_2 & & x_2 & & f_2 \\ w & w & w & & \dots & = & \dots & , \\ a_{n-1} & b_{n-1} & c_{n-1} & x_{n-1} & & f_{n-1} \\ a_n & b_n & x_n & & & f_n \end{array} \quad (4.12)$$

简记为 $\mathbf{Ax} = \mathbf{f}$. 其中, 当 $|i - j| > 1$ 时, $a_{ij} = 0$, 且:

- (a) $|b| > |c| > 0$;
- (b) $|b| = |a_i| + |c_i|$, $a_i, c_i \neq 0$ ($i = 2, 3, \dots, n-1$);
- (c) $|b_n| > |a_n| > 0$.

我们利用矩阵的直接三角分解法来推导解三对角线方程组(4.12)的计算公式. 由系数阵 \mathbf{A} 的特点, 可以将 \mathbf{A} 分解为两个三角阵的乘积, 即

$$\mathbf{A} = \mathbf{LU},$$

其中 L 为下三角矩阵, U 为单位上三角矩阵. 下面我们来说明这种分解是可能的. 设

$$\begin{aligned} \mathbf{A} &= \begin{matrix} b & c \\ a & b & c \\ & w & w & w \\ & a_{n-1} & b_{n-1} & c_{n-1} \\ & a_n & b_n & \\ & 1 & & 1 & \\ & r_2 & r_2 & & 1 & w \\ & w & w & & w & w \\ & r_n & r_n & & & 1 \end{matrix}, \\ &= \begin{matrix} 1 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & 1 & \\ & & & & 1 \end{matrix}, \end{aligned} \quad (4.13)$$

其中 r_i, a_i, r_n 为待定系数. 比较 (4.13) 两边即得

$$\begin{aligned} b_i &= r_i, a_i = r_{i-1}, \\ a_i &= r_i, b_i = r_{i-1} + r_i \quad (i = 2, \dots, n), \\ a &= r_{i-1} \quad (i = 2, 3, \dots, n-1). \end{aligned} \quad (4.14)$$

由 $r_1 = b_1 - 0, |b_1| > |c_1| > 0, r_1 = a_1/b_1$, 得 $0 < |r_1| < 1$. 下面我们用归纳法证明

$$|r_i| > |a_i| > 0 \quad (i = 1, 2, \dots, n-1), \quad (4.15)$$

即 $0 < |r_i| < 1$, 从而由 (4.14) 可求出 r_i .

(4.15) 对 $i=1$ 是成立的. 现设 (4.15) 对 $i-1$ 成立, 求证对 i 亦成立.

由归纳法假设 $0 < |r_{i-1}| < 1$, 又由 (4.15) 及 \mathbf{A} 的假设条件有

$$\begin{aligned} |r_i| &= |b_i - a_{i-1}| / |b_{i-1}| > |b_i| / |a_{i-1}| \\ &> |b_i| / |a_i| / |c_i| = 0, \end{aligned}$$

也就是 $0 < |r_i| < 1$. 由 (4.14) 得到

$$r_i = b_i - a_{i-1} \quad (i = 2, \dots, n);$$

$$r_i = c / (b_i - a_{i-1}) \quad (i = 2, 3, \dots, n-1).$$

这就是说,由 \mathbf{A} 的假设条件,我们完全确定了 $\{a_i\}$, $\{r_i\}$, $\{c_i\}$, 实现了 \mathbf{A} 的 LU 分解.

求解 $\mathbf{Ax} = \mathbf{f}$ 等价于解两个三角形方程组

(1) $\mathbf{Ly} = \mathbf{f}$, 求 \mathbf{y} ; (2) $\mathbf{Ux} = \mathbf{y}$, 求 \mathbf{x} .

从而得到解三对角线方程组的追赶法公式:

1. 计算 $\{r_i\}$ 的递推公式

$$r_1 = c / b_1,$$

$$r_i = c / (b_i - a_{i-1}) \quad (i = 2, 3, \dots, n-1);$$

2. 解 $\mathbf{Ly} = \mathbf{f}$

$$y_1 = f_1 / b_1,$$

$$y_i = (f_i - a_{i-1}y_{i-1}) / (b_i - a_{i-1}) \quad (i = 2, 3, \dots, n);$$

3. 解 $\mathbf{Ux} = \mathbf{y}$

$$x_n = y_n,$$

$$x_i = y_i - r_i x_{i+1} \quad (i = n-1, n-2, \dots, 2, 1).$$

我们将计算系数 r_1, r_2, \dots, r_{n-1} 及 y_1, y_2, \dots, y_n 的过程称为追的过程, 将计算方程组的解 x_n, x_{n-1}, \dots, x_1 的过程称为赶的过程.

总结上述讨论有

定理 12 设有三对角线方程组 $\mathbf{Ax} = \mathbf{f}$, 其中 \mathbf{A} 满足条件

(a), (b), (c), 则 \mathbf{A} 为非奇异矩阵且追赶法计算公式中的 $\{r_i\}$, $\{c_i\}$ 满足:

$$1^\circ 0 < |r_i| < 1 \quad (i = 1, 2, \dots, n-1);$$

$$2^\circ 0 < |c_i| - |b_i| - |a_i| < |r_i| < |b_i| + |a_i| \quad (i = 2, 3, \dots, n-1);$$

$$0 < |b_n| - |a_n| < |r_n| < |b_n| + |a_n|.$$

追赶法公式实际上就是把高斯消去法用到求解三对角线方程

组上去的结果. 这时由于 \mathbf{A} 特别简单, 因此使得求解的计算公式非常简单, 而且计算量仅为 $5n - 4$ 次乘除法, 而另外增加解一个方程组 $\mathbf{Ax} = \mathbf{f}$ 仅增加 $3n - 2$ 次乘除运算. 易见追赶法的计算量是比较小的.

由定理 12 的 1°, 2° 说明追赶法计算公式中不会出现中间结果数量级的巨大增长和舍入误差的严重累积.

在计算机实现时我们只需用三个一维数组分别存储 \mathbf{A} 的三条线元素 $\{a_i\}$, $\{b\}$, $\{c\}$, 此外还需要用两组工作单元保存 $\{x_i\}$, $\{y_i\}$ 或 $\{z_i\}$.

5.5 向量和矩阵的范数

为了研究线性方程组近似解的误差估计和迭代法的收敛性, 我们需要对 \mathbf{R}^n (n 维向量空间) 中向量(或 $\mathbf{R}^{n \times n}$ 中矩阵)的“大小”引进某种度量——向量(或矩阵)范数的概念. 向量范数概念是三维欧氏空间中向量长度概念的推广, 在数值分析中起着重要作用.

首先将向量长度概念推广到 \mathbf{R}^n (或 \mathbf{C}^n) 中.

定义 1 设 $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$, $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$
 \mathbf{R}^n (或 \mathbf{C}^n). 将实数 $(\mathbf{x}, \mathbf{y}) = \mathbf{y}^T \mathbf{x} = \sum_{i=1}^n x_i y_i$ (或复数 $(\mathbf{x}, \mathbf{y}) = \mathbf{y}^H \mathbf{x} = \sum_{i=1}^n x_i \bar{y}_i$) 称为向量 \mathbf{x}, \mathbf{y} 的数量积.

将非负实数 $\|\mathbf{x}\|_2 = (\mathbf{x}, \mathbf{x})^{\frac{1}{2}} = (\sum_{i=1}^n x_i^2)^{\frac{1}{2}}$ 称为向量 \mathbf{x} 的欧氏范数.

下述定理可在线性代数书中找到.

定理 13 设 $\mathbf{x}, \mathbf{y} \in \mathbf{R}^n$ (或 \mathbf{C}^n), 则

1. $(\mathbf{x}, \mathbf{x}) = 0$, 当且仅当 $\mathbf{x} = \mathbf{0}$ 时成立;

2. $(\mathbf{x}, \mathbf{y}) = (\mathbf{x}, \mathbf{y})$. 为实数 (或 $(\mathbf{x}, \mathbf{y}) = \bar{(\mathbf{x}, \mathbf{y})}$, 为复数);

3. $(\mathbf{x}, \mathbf{y}) = (\mathbf{y}, \mathbf{x})$ (或 $(\mathbf{x}, \mathbf{y}) = (\overline{\mathbf{y}}, \mathbf{x})$);

4. $(\mathbf{x}_1 + \mathbf{x}_2, \mathbf{y}) = (\mathbf{x}_1, \mathbf{y}) + (\mathbf{x}_2, \mathbf{y})$;

5. (Cauchy-Schwarz 不等式)

$$|(\mathbf{x}, \mathbf{y})| \leq \|\mathbf{x}\|_2 \cdot \|\mathbf{y}\|_2,$$

等式当且仅当 \mathbf{x} 与 \mathbf{y} 线性相关时成立;

6. 三角不等式

$$\|\mathbf{x} + \mathbf{y}\|_2 \leq \|\mathbf{x}\|_2 + \|\mathbf{y}\|_2.$$

我们还可以用其他办法来度量 \mathbf{R}^n 中向量的“大小”. 例如, 对于 $\mathbf{x} = (x_1, x_2)^T \in \mathbf{R}^2$, 可以用一个 \mathbf{x} 的函数 $N(\mathbf{x}) = \max_{i=1,2} |x_i|$ 来度量 \mathbf{x} 的“大小”, 而且这种度量 \mathbf{x} 大小”的方法计算起来比欧氏范数方便. 在许多应用中, 对度量向量 \mathbf{x} 大小”的函数 $N(\mathbf{x})$ 都要求是正定的、齐次的且满足三角不等式. 下面我们给出向量范数的一般定义.

定义 2(向量的范数) 如果向量 $\mathbf{x} \in \mathbf{R}^n$ (或 \mathbf{C}^n) 的某个实值函数 $N(\mathbf{x}) = \|\mathbf{x}\|$, 满足条件:

(1) $\mathbf{x} = 0$ ($\mathbf{x} = 0$ 当且仅当 $\mathbf{x} = \mathbf{0}$) (正定条件),

(2) $\mathbf{x} = |\mathbf{x}|$ ($\mathbf{x} \in \mathbf{R}$ 或 \mathbf{C}), (5.1)

(3) $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ (三角不等式),

则称 $N(\mathbf{x})$ 是 \mathbf{R}^n (或 \mathbf{C}^n) 上的一个向量范数(或模). 由(3)可推出不等式

(4) $|\mathbf{x} - \mathbf{y}| \geq \|\mathbf{x} - \mathbf{y}\|$. (5.2)

下面我们给出几种常用的向量范数.

1. 向量的 ∞ -范数(最大范数):

$$\|\mathbf{x}\|_1 = \max_{i=1}^n |x_i|.$$

容易验证这样定义的向量 \mathbf{x} 的函数 $N(\mathbf{x}) = \|\mathbf{x}\|_1$ 满足向量范数的三个条件 .

2. 向量的 1 - 范数:

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

同样可证 $N(\mathbf{x}) = \|\mathbf{x}\|_1$ 是 \mathbb{R}^n 上的一个向量范数 .

3. 向量的 2 - 范数:

$$\|\mathbf{x}\|_2 = (\mathbf{x} \cdot \mathbf{x})^{1/2} = \left(\sum_{i=1}^n x_i^2 \right)^{1/2}.$$

由定理 13 知 $N(\mathbf{x}) = \|\mathbf{x}\|_2$ 是 \mathbb{R}^n 上一个向量范数, 称为向量 \mathbf{x} 的欧氏范数 .

4. 向量的 p - 范数:

$$\|\mathbf{x}\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p},$$

其中 $p \in [1, \infty)$, 可以证明向量函数 $N(\mathbf{x}) = \|\mathbf{x}\|_p$ 是 \mathbb{R}^n 上向量的范数且容易说明上述三种范数是 p - 范数的特殊情况 ($\|\mathbf{x}\|_1 =$

$$\lim_{p \rightarrow 1} \|\mathbf{x}\|_p).$$

例 6 计算向量 $\mathbf{x} = (1, -2, 3)^T$ 的各种范数 .

解 $\|\mathbf{x}\|_1 = 6$, $\|\mathbf{x}\|_2 = 3$, $\|\mathbf{x}\|_\infty = 14$.

定义 3 设 $\{\mathbf{x}^{(k)}\}$ 为 \mathbb{R}^n 中一向量序列, $\mathbf{x}^* \in \mathbb{R}^n$, 记 $\mathbf{x}^{(k)} = (x_1^{(k)}, x_2^{(k)}, \dots, x_n^{(k)})^T$, $\mathbf{x}^* = (x_1^*, x_2^*, \dots, x_n^*)^T$. 如果 $\lim_k x_i^{(k)} = x_i^*$ ($i = 1, 2, \dots, n$), 则称 $\mathbf{x}^{(k)}$ 收敛于向量 \mathbf{x}^* , 记为

$$\lim_k \mathbf{x}^{(k)} = \mathbf{x}^*.$$

定理 14 ($N(\mathbf{x})$ 的连续性) 设非负函数 $N(\mathbf{x}) = \|\mathbf{x}\|$ 为 \mathbb{R}^n 上任一向量范数, 则 $N(\mathbf{x})$ 是 \mathbf{x} 的分量 x_1, x_2, \dots, x_n 的连续函数 .

证明 设 $\mathbf{x} = \sum_{i=1}^n x_i \mathbf{e}_i$, $\mathbf{y} = \sum_{i=1}^n y_i \mathbf{e}_i$, 其中 $\mathbf{e} = (0, \dots, 1, 0, \dots, 0)^T$.

只须证明当 $\mathbf{x} \parallel \mathbf{y}$ 时 $N(\mathbf{x}) = N(\mathbf{y})$ 即成. 事实上

$$\begin{aligned} |N(\mathbf{x}) - N(\mathbf{y})| &= |\mathbf{x} - \mathbf{y}| / \|\mathbf{x} - \mathbf{y}\| \\ &= \left(\sum_{i=1}^n (x_i - y_i) \right) \mathbf{e} \\ &= \left\| \sum_{i=1}^n (x_i - y_i) \mathbf{e} \right\|, \end{aligned}$$

即 $|N(\mathbf{x}) - N(\mathbf{y})| = c \|\mathbf{x} - \mathbf{y}\| = 0$ (当 $\mathbf{x} \parallel \mathbf{y}$ 时),

其中 $c = \left\| \sum_{i=1}^n \mathbf{e}_i \right\|$.

定理 15(向量范数的等价性) 设 $\mathbf{x}_s, \mathbf{x}_t$ 为 \mathbf{R}^n 上向量的任意两种范数, 则存在常数 $a, c > 0$, 使得对一切 $\mathbf{x} \in \mathbf{R}^n$ 有

$$a \mathbf{x}_s \leq \mathbf{x}_t \leq c \mathbf{x}_s.$$

证明 只要就 $\mathbf{x}_s = \mathbf{x}$ 证明上式成立即可, 即证明存在常数 $a, c > 0$, 使

$$a \frac{\mathbf{x}_t}{\mathbf{x}} \leq c, \text{ 对一切 } \mathbf{x} \in \mathbf{R}^n \text{ 且 } \mathbf{x} \neq \mathbf{0}.$$

考虑泛函

$$f(\mathbf{x}) = \mathbf{x}_t / \|\mathbf{x}\|, \quad \mathbf{x} \in \mathbf{R}^n.$$

记 $S = \{\mathbf{x} \mid \mathbf{x}_s = 1, \mathbf{x} \in \mathbf{R}^n\}$, 则 S 是一个有界闭集. 由于 $f(\mathbf{x})$ 为 S 上的连续函数, 所以 $f(\mathbf{x})$ 于 S 上达到最大最小值, 即存在 $\mathbf{x}, \mathbf{x} \in S$ 使得

$$f(\mathbf{x}) = \min_{\mathbf{x} \in S} f(\mathbf{x}) = a, \quad f(\mathbf{x}) = \max_{\mathbf{x} \in S} f(\mathbf{x}) = c.$$

设 $\mathbf{x} \in \mathbf{R}^n$ 且 $\mathbf{x} \neq \mathbf{0}$, 则 $\frac{\|\mathbf{x}\|}{\|\mathbf{x}\|_t} = S$, 从而有

$$a = f\left(\frac{\mathbf{x}}{\|\mathbf{x}\|_t}\right) < c, \quad (5.3)$$

显然 $a, c > 0$, 上式为

$$a = \left\| \frac{\mathbf{x}}{\|\mathbf{x}\|_t} \right\|_t < c,$$

即

$$a \|\mathbf{x}\|_t < c \|\mathbf{x}\|, \text{ 对一切 } \mathbf{x} \in \mathbf{R}^n.$$

注意, 定理 15 不能推广到无穷维空间. 由定理 15 可得到结论: 如果在一种范数意义下向量序列收敛时, 则在任何一种范数意义下该向量序列均收敛.

定理 16 $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^*$ $\lim_{k \rightarrow \infty} \|\mathbf{x}^{(k)} - \mathbf{x}^*\| = 0$, 其中 \cdot 为向量的任一种范数.

证明 显然, $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^*$ $\lim_{k \rightarrow \infty} \|\mathbf{x}^{(k)} - \mathbf{x}^*\| = 0$, 而对于 \mathbf{R}^n 上任一种范数 \cdot , 由定理 15, 存在常数 $a, c > 0$ 使

$$a \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_t < c \|\mathbf{x}^{(k)} - \mathbf{x}^*\|, \quad ,$$

于是又有

$$\lim_{k \rightarrow \infty} \|\mathbf{x}^{(k)} - \mathbf{x}^*\| = 0 \quad \lim_{k \rightarrow \infty} \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_t = 0.$$

下面我们将向量范数概念推广到矩阵上去. 视 $\mathbf{R}^{n \times n}$ 中的矩阵为 \mathbf{R}^{n^2} 中的向量, 则由 \mathbf{R}^{n^2} 上的 2-范数可以得到 $\mathbf{R}^{n \times n}$ 中矩阵的一种范数

$$F(\mathbf{A}) = \|\mathbf{A}\|_F = \sqrt{\sum_{i,j=1}^n a_{i,j}^2},$$

称为 \mathbf{A} 的 Frobenius 范数. $\|\mathbf{A}\|_F$ 显然满足正定性、齐次性及三角不等式.

下面我们给出矩阵范数的一般定义.

定义 4(矩阵的范数) 如果矩阵 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 的某个非负的实值函数 $N(\mathbf{A}) = \|\mathbf{A}\|$, 满足条件

- (1) $\mathbf{A} \geq 0$ ($\mathbf{A} = 0 \iff \mathbf{A} = \mathbf{0}$) (正定条件);
- (2) $c\mathbf{A} = |c| \|\mathbf{A}\|$, c 为实数 (齐次条件); (5.4)
- (3) $\mathbf{A} + \mathbf{B} \leq \|\mathbf{A}\| + \|\mathbf{B}\|$ (三角不等式);
- (4) $\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|$.

则称 $N(\mathbf{A})$ 是 $\mathbb{R}^{n \times n}$ 上的一个矩阵范数(或模).

上面我们定义的 $F(\mathbf{A}) = \|\mathbf{A}\|_F$ 就是 $\mathbb{R}^{n \times n}$ 上的一个矩阵范数.

由于在大多数与估计有关的问题中, 矩阵和向量会同时参与讨论, 所以希望引进一种矩阵的范数, 它是和向量范数相联系而且和向量范数相容的, 即对任何向量 $\mathbf{x} \in \mathbb{R}^n$ 及 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 都成立

$$\|\mathbf{Ax}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|. \quad (5.5)$$

为此我们再引进一种矩阵的范数.

定义 5(矩阵的算子范数) 设 $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{A} \in \mathbb{R}^{n \times n}$, 给出一种向量范数 $\|\mathbf{x}\|_v$ (如 $v=1, 2$ 或 ∞), 相应地定义一个矩阵的非负函数

$$\|\mathbf{A}\|_v = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|_v}{\|\mathbf{x}\|_v}. \quad (5.6)$$

可验证 $\|\mathbf{A}\|_v$ 满足定义 4(见下面定理), 所以 $\|\mathbf{A}\|_v$ 是 $\mathbb{R}^{n \times n}$ 上矩阵的一个范数, 称为 \mathbf{A} 的算子范数.

定理 17 设 $\|\mathbf{x}\|_v$ 是 \mathbb{R}^n 上一个向量范数, 则 $\|\mathbf{A}\|_v$ 是 $\mathbb{R}^{n \times n}$ 上矩阵的范数, 且满足相容条件

$$\|\mathbf{Ax}\|_v \leq \|\mathbf{A}\|_v \|\mathbf{x}\|_v. \quad (5.7)$$

证明 由(5.6)相容性条件(5.7)是显然的. 现只验证定义 4 中条件(4).

由(5.7), 有

$$\|\mathbf{ABx}\|_v \leq \|\mathbf{A}\|_v \|\mathbf{Bx}\|_v \leq \|\mathbf{A}\|_v \|\mathbf{B}\|_v \|\mathbf{x}\|_v.$$

当 $\mathbf{x} = \mathbf{0}$ 时, 有

$$\frac{\|\mathbf{ABx}\|_v}{\|\mathbf{x}\|_v} = \|\mathbf{A}\|_v \|\mathbf{B}\|_v,$$

故

$$\|\mathbf{AB}\|_v = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{ABx}\|_v}{\|\mathbf{x}\|_v} = \|\mathbf{A}\|_v \|\mathbf{B}\|_v.$$

显然这种矩阵的范数 $\|\mathbf{A}\|_v$ 依赖于向量范数 $\|\mathbf{x}\|_v$ 的具体含义. 也就是说, 当给出一种具体的向量范数 $\|\mathbf{x}\|_v$ 时, 相应地就得到了一种矩阵范数 $\|\mathbf{A}\|_v$.

定理 18 设 $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{A} \in \mathbb{R}^{n \times n}$, 则

1. $\|\mathbf{A}\|_1 = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|$ (称为 \mathbf{A} 的行范数),
2. $\|\mathbf{A}\|_\infty = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}|$ (称为 \mathbf{A} 的列范数),
3. $\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A}^T \mathbf{A})}$ (称为 \mathbf{A} 的 2 - 范数),

其中 $\lambda_{\max}(\mathbf{A}^T \mathbf{A})$ 表示 $\mathbf{A}^T \mathbf{A}$ 的最大特征值.

证明 只就 1, 3 给出证明, 2 同理.

1. 设 $\mathbf{x} = (x_1, \dots, x_n)^T \neq \mathbf{0}$, 不妨设 $\mathbf{A} \neq \mathbf{0}$. 记

$$t = \|\mathbf{x}\|_1 = \max_{1 \leq i \leq n} |x_i|, \quad \mu = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|,$$

则

$$\begin{aligned} \|\mathbf{Ax}\|_1 &= \max_{1 \leq i \leq n} \left| \sum_{j=1}^n a_{ij} x_j \right| = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| |x_j| \\ &\leq t \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| = t \mu. \end{aligned}$$

这说明对任何非零 $\mathbf{x} \in \mathbb{R}^n$, 有

$$\frac{\|\mathbf{Ax}\|_1}{\|\mathbf{x}\|_1} \leq \mu. \quad (5.8)$$

下面来说明有一向量 $\mathbf{x} \neq \mathbf{0}$, 使 $\frac{\|\mathbf{Ax}\|}{\|\mathbf{x}\|} = \mu$. 设 $\mu =$

n
 $\sum_{j=1}^n |a_{0j}|$, 取向量 $\mathbf{x} = (x_1, \dots, x_n)^T$, 其中 $x_j = \operatorname{sgn}(a_{0j})$
 $(j = 1, 2, \dots, n)$. 显然 $\|\mathbf{x}\| = 1$, 且 \mathbf{Ax} 的第 i 个分量为

$$\sum_{i=1}^n a_{ij} x_i = \sum_{j=1}^n |a_{ij}|, \text{ 这说明}$$

$$\|\mathbf{Ax}\| = \max_{1 \leq i \leq n} \left| \sum_{j=1}^n a_{ij} x_j \right| = \sum_{j=1}^n |a_{0j}| = \mu.$$

3. 由于对一切 $\mathbf{x} \in \mathbb{R}^n$, $\|\mathbf{Ax}\|_2^2 = (\mathbf{Ax}, \mathbf{Ax}) = (\mathbf{A}^T \mathbf{Ax}, \mathbf{x}) = 0$,
 从而 $\mathbf{A}^T \mathbf{A}$ 的特征值为非负实数, 设为

$$1, 2, \dots, n, 0. \quad (5.9)$$

$\mathbf{A}^T \mathbf{A}$ 为对称矩阵, 设 $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n$ 为 $\mathbf{A}^T \mathbf{A}$ 的相应于(5.9)的特征向量且 $(\mathbf{u}_i, \mathbf{u}_j) = \delta_{ij}$, 又设 $\mathbf{x} \in \mathbb{R}^n$ 为任一非零向量, 于是有

$$\mathbf{x} = \sum_{i=1}^n c_i \mathbf{u}_i,$$

其中 c_i 为组合系数, 则

$$\frac{\|\mathbf{Ax}\|_2^2}{\|\mathbf{x}\|_2^2} = \frac{(\mathbf{A}^T \mathbf{Ax}, \mathbf{x})}{(\mathbf{x}, \mathbf{x})} = \frac{\sum_{i=1}^n c_i^2 \lambda_i}{\sum_{i=1}^n c_i^2} = \lambda_1.$$

另一方面, 取 $\mathbf{x} = \mathbf{u}_1$, 则上式等号成立, 故

$$\|\mathbf{A}\|_2 = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|_2^2}{\|\mathbf{x}\|_2^2} = \lambda_1 = \max(\mathbf{A}^T \mathbf{A}).$$

由定理 18 看出, 计算一个矩阵的 $\|\mathbf{A}\|_1$, $\|\mathbf{A}\|_\infty$ 还是比较容易的, 而矩阵的 2-范数 $\|\mathbf{A}\|_2$ 在计算上不方便, 但是矩阵的 2-范数具有许多好的性质, 它在理论上是非常有用的.

例 7 设 $\mathbf{A} = \begin{pmatrix} 1 & -2 \\ -3 & 4 \end{pmatrix}$, 计算 \mathbf{A} 的各种范数.

$$\text{解 } \quad \mathbf{A}_1 = 6, \quad \mathbf{A} = 7, \quad \mathbf{A}_F = 5.477,$$

$$\mathbf{A}_2 = 15 + 221 \cdot 5.46.$$

我们指出, 对于复矩阵(即 $\mathbf{A} \in \mathbb{C}^{n \times n}$)定理 18 中 1, 2. 显然也成立, 对于 3 应改为

$$\mathbf{A}_2 = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^H \mathbf{A}^H \mathbf{A} \mathbf{x}}{\mathbf{x}^H \mathbf{x}}^{1/2} = \max_{\mathbf{x} \neq \mathbf{0}} (\mathbf{A}^H \mathbf{A}).$$

定义 6 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 的特征值为 λ_i ($i = 1, 2, \dots, n$), 称

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

为 \mathbf{A} 的谱半径.

定理 19(特征值上界) 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$, 则 $(\mathbf{A}) \leq \mathbf{A}_F$, 即 \mathbf{A} 的谱半径不超过 \mathbf{A} 的任何一种算子范数(对 \mathbf{A}_F 亦对).

证明 设 λ 是 \mathbf{A} 的任一特征值, \mathbf{x} 为相应的特征向量, 则 $\mathbf{Ax} = \lambda \mathbf{x}$, 由(5.7)得

$$|\lambda| / \|\mathbf{x}\| = \|\mathbf{x}\| = \|\mathbf{Ax}\| / \|\mathbf{A}\| \|\mathbf{x}\|,$$

注意到 $\mathbf{x} \neq \mathbf{0}$, 即得

$$|\lambda| / \|\mathbf{A}\|.$$

定理 20 如果 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 为对称矩阵, 则 $\mathbf{A}_2 = (\mathbf{A})$.

证明留作习题.

定理 21 如果 $\mathbf{B} < 1$, 则 $\mathbf{I} \pm \mathbf{B}$ 为非奇异矩阵, 且

$$(\mathbf{I} \pm \mathbf{B})^{-1} = \frac{1}{1 - \|\mathbf{B}\|},$$

其中 $\|\cdot\|$ 是指矩阵的算子范数.

证明 用反证法. 若 $\det(\mathbf{I} - \mathbf{B}) = 0$, 则 $(\mathbf{I} - \mathbf{B})\mathbf{x} = \mathbf{0}$ 有非零解, 即存在 $\mathbf{x} \neq \mathbf{0}$ 使 $\mathbf{Bx} = \mathbf{x}$, $\frac{\mathbf{Bx}}{\|\mathbf{x}\|} = 1$, 故 $\|\mathbf{B}\| = 1$, 与假设矛盾. 又由 $(\mathbf{I} - \mathbf{B})(\mathbf{I} - \mathbf{B})^{-1} = \mathbf{I}$, 有

$$(\mathbf{I} - \mathbf{B})^{-1} = \mathbf{I} + \mathbf{B}(\mathbf{I} - \mathbf{B})^{-1},$$

从而

$$\begin{aligned} (\mathbf{I} \pm \mathbf{B})^{-1} &= \mathbf{I} - \frac{\mathbf{B}}{1 + \frac{1}{\mathbf{B}}} , \\ (\mathbf{I} - \mathbf{B})^{-1} &= \frac{1}{1 - \frac{1}{\mathbf{B}}} . \end{aligned}$$

5.6 误差分析

5.6.1 矩阵的条件数

考虑线性方程组

$$\mathbf{Ax} = \mathbf{b},$$

其中设 \mathbf{A} 为非奇异矩阵, \mathbf{x} 为方程组的精确解.

由于 \mathbf{A} (或 \mathbf{b}) 元素是测量得到的, 或者是计算的结果, 在第一种情况 \mathbf{A} (或 \mathbf{b}) 常带有某些观测误差, 在后一种情况 \mathbf{A} (或 \mathbf{b}) 又包含有舍入误差. 因此我们处理的实际矩阵是 $\mathbf{A} + \Delta \mathbf{A}$ (或 $\mathbf{b} + \Delta \mathbf{b}$), 下面我们来研究数据 \mathbf{A} (或 \mathbf{b}) 的微小误差对解的影响. 即考虑估计 $\mathbf{x} - \mathbf{y}$, 其中 \mathbf{y} 是 $(\mathbf{A} + \Delta \mathbf{A})\mathbf{y} = \mathbf{b}$ 的解.

首先考察一个例子.

例 8 设有方程组

$$\begin{array}{cc} 1 & 1 \\ 1 & 1 \end{array} \begin{array}{c} x_1 \\ x_2 \end{array} = \begin{array}{c} 2 \\ 2 \end{array}. \quad (6.1)$$

记为 $\mathbf{Ax} = \mathbf{b}$, 它的精确解为 $\mathbf{x} = (2, 0)^T$.

现在考虑常数项的微小变化对方程组解的影响, 即考察方程组

$$\begin{array}{cc} 1 & 1 \\ 1 & 1 \end{array} \begin{array}{c} y_1 \\ y_2 \end{array} = \begin{array}{c} 2 \\ 2.0001 \end{array}, \quad (6.2)$$

也可表示为 $\mathbf{A}(\mathbf{x} + \Delta \mathbf{x}) = \mathbf{b} + \Delta \mathbf{b}$, 其中 $\mathbf{b} = (0, 2.0001)^T$, $\Delta \mathbf{b} = (0, 0.0001)^T$, $\Delta \mathbf{x} = \mathbf{x} + \mathbf{x}'$, \mathbf{x}' 为(6.1)的解. 显然方程组(6.2)的解为 $\mathbf{x} + \Delta \mathbf{x} = (1, 1)^T$.

我们看到(6.1)的常数项 \mathbf{b} 的第 2 个分量只有 $\frac{1}{10000}$ 的微小变化, 方程组的解却变化很大. 这样的方程组称为病态方程组.

定义 7 如果矩阵 \mathbf{A} 或常数项 \mathbf{b} 的微小变化, 引起方程组 $\mathbf{Ax} = \mathbf{b}$ 解的巨大变化, 则称此方程组为“病态”方程组, 矩阵 \mathbf{A} 称为“病态”矩阵(相对于方程组而言), 否则称方程组为“良态”方程组, \mathbf{A} 称为“良态”矩阵.

应该注意, 矩阵的“病态”性质是矩阵本身的特性, 下面我们希望找出刻画矩阵“病态”性质的量. 设有方程组

$$\mathbf{Ax} = \mathbf{b}, \quad (6.3)$$

其中 \mathbf{A} 为非奇异阵, \mathbf{x} 为(6.3)的准确解. 以下我们研究方程组的系数矩阵 \mathbf{A} (或 \mathbf{b})的微小误差(扰动)时对解的影响.

现设 \mathbf{A} 是精确的, \mathbf{b} 有误差 \mathbf{b}' , 解为 \mathbf{x}' , 则

$$\begin{aligned} \mathbf{A}(\mathbf{x}' + \mathbf{x}) &= \mathbf{b}' + \mathbf{b}, \quad \mathbf{x}' = \mathbf{A}^{-1} \mathbf{b}', \\ \mathbf{x} &\qquad \mathbf{A}^{-1} \qquad \mathbf{b}' . \end{aligned} \quad (6.4)$$

由(6.3)

$$\begin{aligned} \mathbf{b} &\qquad \mathbf{A} \qquad \mathbf{x} , \\ \frac{1}{\mathbf{x}} &\qquad -\frac{\mathbf{A}}{\mathbf{b}} \quad (\text{设 } \mathbf{b} = \mathbf{0}) . \end{aligned} \quad (6.5)$$

于是由(6.4)及(6.5), 得

定理 22 设 \mathbf{A} 是非奇异阵, $\mathbf{Ax} = \mathbf{b} = \mathbf{0}$, 且

$$\mathbf{A}(\mathbf{x}' + \mathbf{x}) = \mathbf{b}' + \mathbf{b},$$

则 $\frac{\mathbf{x}'}{\mathbf{x}} = \mathbf{A}^{-1} \mathbf{A} = \frac{\mathbf{b}'}{\mathbf{b}} .$

上式给出了解的相对误差的上界, 常数项 \mathbf{b} 的相对误差在解中可能放大 $\mathbf{A}^{-1} \mathbf{A}$ 倍.

现设 \mathbf{b} 是精确的, \mathbf{A} 有微小误差(扰动) \mathbf{A}' , 解为 $\mathbf{x}' + \mathbf{x}$, 则

$$(\mathbf{A}' + \mathbf{A})(\mathbf{x}' + \mathbf{x}) = \mathbf{b},$$

$$(\mathbf{A} + \mathbf{A}) \mathbf{x} = -(\mathbf{A}) \mathbf{x}. \quad (6.6)$$

如果 \mathbf{A} 不受限制的话, $\mathbf{A} + \mathbf{A}$ 可能奇异, 而

$$(\mathbf{A} + \mathbf{A}) = \mathbf{A}(\mathbf{I} + \mathbf{A}^{-1} \mathbf{A}),$$

由定理 21 知, 当 $\|\mathbf{A}^{-1} \mathbf{A}\| < 1$ 时, $(\mathbf{I} + \mathbf{A}^{-1} \mathbf{A})^{-1}$ 存在. 由(6.6)式

$$\mathbf{x} = -(\mathbf{I} + \mathbf{A}^{-1} \mathbf{A})^{-1} \mathbf{A}^{-1} (\mathbf{A}) \mathbf{x},$$

因此

$$\mathbf{x} = \frac{\mathbf{A}^{-1} \mathbf{A} \mathbf{x}}{1 - \mathbf{A}^{-1} (\mathbf{A})}.$$

设 $\|\mathbf{A}^{-1} \mathbf{A}\| < 1$, 即得

$$\frac{\mathbf{x}}{\mathbf{x}} = \frac{\mathbf{A}^{-1} \mathbf{A} \frac{\mathbf{A}}{\mathbf{A}}}{1 - \mathbf{A}^{-1} \mathbf{A} \frac{\mathbf{A}}{\mathbf{A}}}. \quad (6.7)$$

定理 23 设 \mathbf{A} 为非奇异矩阵, $\mathbf{A}\mathbf{x} = \mathbf{b}$, 且

$$(\mathbf{A} + \mathbf{A})(\mathbf{x} + \mathbf{x}) = \mathbf{b}.$$

如果 $\|\mathbf{A}^{-1} \mathbf{A}\| < 1$, 则(6.7)式成立.

如果 \mathbf{A} 充分小, 且在条件 $\|\mathbf{A}^{-1} \mathbf{A}\| < 1$ 下, 那么(6.7)式

说明矩阵 \mathbf{A} 的相对误差 $\frac{\mathbf{A}}{\mathbf{A}}$ 在解中可能放大 $\|\mathbf{A}^{-1} \mathbf{A}\|$ 倍.

总之, 量 $\|\mathbf{A}^{-1} \mathbf{A}\|$ 愈小, 由 \mathbf{A} (或 \mathbf{b}) 的相对误差引起的解的相对误差就愈小; 量 $\|\mathbf{A}^{-1} \mathbf{A}\|$ 愈大, 解的相对误差就可能愈大. 所以量 $\|\mathbf{A}^{-1} \mathbf{A}\|$ 实际上刻画了解对原始数据变化的灵敏程度, 即刻画了方程组的“病态”程度, 于是引进下述定义:

定义 8 设 \mathbf{A} 为非奇异阵, 称数 $\text{cond}(\mathbf{A})_v = \|\mathbf{A}^{-1}\|_v \|\mathbf{A}\|_v$ ($v = 1, 2$ 或 ∞) 为矩阵 \mathbf{A} 的条件数.

由此看出矩阵的条件数与范数有关.

矩阵的条件数是一个十分重要的概念, 由上面讨论知, 当 \mathbf{A} 的条件数相对的大, 即 $\text{cond}(\mathbf{A}) \gg 1$ 时, 则(6.3)是“病态”的(即 \mathbf{A} 是

“病态”矩阵,或者说 \mathbf{A} 是坏条件的,相对于解方程组),当 \mathbf{A} 的条件数相对的小,则(6.3)是“良态”的(或者说 \mathbf{A} 是好条件的).注意,方程组病态性质是方程组本身的特性. \mathbf{A} 的条件数愈大,方程组的病态程度愈严重,也就愈难用一般的计算方法求得比较准确的解.

通常使用的条件数,有

$$(1) \operatorname{cond}(\mathbf{A}) = \|\mathbf{A}^{-1}\| \|\mathbf{A}\|;$$

(2) \mathbf{A} 的谱条件数

$$\operatorname{cond}(\mathbf{A})_2 = \|\mathbf{A}\|_2 \|\mathbf{A}^{-1}\|_2 = \frac{\max_{1 \leq i \leq n} (\mathbf{A}^T \mathbf{A})_{ii}}{\min_{1 \leq i \leq n} (\mathbf{A}^T \mathbf{A})_{ii}}.$$

当 \mathbf{A} 为对称矩阵时

$$\operatorname{cond}(\mathbf{A})_2 = \frac{| \lambda_1 |}{| \lambda_n |},$$

其中 λ_1, λ_n 为 \mathbf{A} 的绝对值最大和绝对值最小的特征值.

条件数的性质:

1. 对任何非奇异矩阵 \mathbf{A} 都有 $\operatorname{cond}(\mathbf{A})_v \geq 1$. 事实上,

$$\operatorname{cond}(\mathbf{A})_v = \|\mathbf{A}^{-1}\|_v \|\mathbf{A}\|_v = \|\mathbf{A}^{-1}\mathbf{A}\|_v = 1;$$

2. 设 \mathbf{A} 为非奇异阵且 $c > 0$ (常数), 则

$$\operatorname{cond}(c\mathbf{A})_v = \operatorname{cond}(\mathbf{A})_v;$$

3. 如果 \mathbf{A} 为正交矩阵, 则 $\operatorname{cond}(\mathbf{A})_2 = 1$; 如果 \mathbf{A} 为非奇异矩阵, \mathbf{R} 为正交矩阵, 则

$$\operatorname{cond}(\mathbf{R}\mathbf{A})_2 = \operatorname{cond}(\mathbf{A}\mathbf{R})_2 = \operatorname{cond}(\mathbf{A})_2.$$

例 9 已知希尔伯特(Hilbert)矩阵

$$\mathbf{H}_n = \begin{matrix} 1 & \frac{1}{2} & \cdots & \frac{1}{n} \\ \frac{1}{2} & \frac{1}{3} & \cdots & \frac{1}{n+1} \\ \cdots & \cdots & \cdots & \cdots \\ \frac{1}{n} & \frac{1}{1+n} & \cdots & \frac{1}{2n-1} \end{matrix},$$

计算 \mathbf{H} 的条件数 .

解

$$\mathbf{H}_3 = \begin{pmatrix} 1 & \frac{1}{2} & \frac{1}{3} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} \end{pmatrix}, \quad \mathbf{H}_3^{-1} = \begin{pmatrix} 9 & -36 & 30 \\ -36 & 192 & -180 \\ 30 & -180 & 180 \end{pmatrix}.$$

(1) 计算 \mathbf{H} 条件数 $\text{cond}(\mathbf{H}_3)$

$$\mathbf{H}_3 = 11/6, \quad \mathbf{H}_3^{-1} = 408, \text{ 所以 } \text{cond}(\mathbf{H}_3) = 748.$$

同样可计算 $\text{cond}(\mathbf{H}) = 2.9 \times 10^7$, $\text{cond}(\mathbf{H}) = 9.85 \times 10^8$. 当 n 愈大时, \mathbf{H} 矩阵病态愈严重 .

(2) 考虑方程组

$$\mathbf{H}_3 \mathbf{x} = (11/6, 13/12, 47/60)^T = \mathbf{b},$$

设 \mathbf{H} 及 \mathbf{b} 有微小误差(取 3 位有效数字)有

$$\begin{array}{cccccc} 1.00 & 0.500 & 0.333 & x_1 + x_1 & 1.83 \\ 0.500 & 0.333 & 0.250 & x_2 + x_2 & = 1.08, & (6.8) \\ 0.333 & 0.250 & 0.200 & x_3 + x_3 & 0.783 \end{array}$$

简记为 $(\mathbf{H}_3 + \mathbf{H}_3)(\mathbf{x} + \mathbf{x}) = \mathbf{b} + \mathbf{b}$. 方程组 $\mathbf{H}_3 \mathbf{x} = \mathbf{b}$ 与 (6.8) 的精确解分别为: $\mathbf{x} = (1, 1, 1)^T$, $\mathbf{x} + \mathbf{x} = (1.089512538, 0.487967062, 1.491002798)^T$. 于是

$$\mathbf{x} = (0.0895, -0.5120, 0.4910)^T,$$

$$\frac{\mathbf{H}_3}{\mathbf{H}_3} = 0.18 \times 10^{-3} < 0.02\%,$$

$$\frac{\mathbf{b}}{\mathbf{b}} = 0.182\%, \quad \frac{\mathbf{x}}{\mathbf{x}} = 51.2\%.$$

这就是说 \mathbf{H} 与 \mathbf{b} 相对误差不超过 0.3%, 而引起解的相对误差超过 50% .

由上面的讨论, 要判别一个矩阵是否病态需要计算条件数 $\text{cond}(\mathbf{A}) = \|\mathbf{A}^{-1}\| \cdot \|\mathbf{A}\|$, 而计算 \mathbf{A}^{-1} 是比较费劲的, 那么在实际计算中如何发现病态情况呢?

(1) 如果在 \mathbf{A} 的三角约化时(尤其是用主元素消去法解(6.3)时)出现小主元, 对大多数矩阵来说, \mathbf{A} 是病态矩阵. 例如用选主元的直接三角分解法解方程组(6.8)(结果舍入为 3 位浮点数), 则有

$$\begin{array}{rcl} & 1 & \\ \mathbf{L}_3 (\mathbf{H}_3 + \mathbf{H}_3) = & 0.333 & 1 \\ & 0.500 & 0.994 & 1 \\ & 1 & 0.5000 & 0.3330 \\ \times & 0.0835 & 0.0891 & = \mathbf{LU} \\ & -0.00507 & & \end{array}$$

(2) 系数矩阵的行列式值相对说很小, 或系数矩阵某些行近似线性相关, 这时 \mathbf{A} 可能病态.

(3) 系数矩阵 \mathbf{A} 元素间数量级相差很大, 并且无一定规则, \mathbf{A} 可能病态.

用选主元素的消去法不能解决病态问题, 对于病态方程组可采用高精度的算术运算(采用双倍字长进行运算)或者采用预处理方法. 即将求解 $\mathbf{Ax} = \mathbf{b}$ 转化为一等价方程组

$$\mathbf{PAQy} = \mathbf{Pb},$$

$$\mathbf{y} = \mathbf{Q}^{-1} \mathbf{x}.$$

选择非奇异矩阵 \mathbf{P}, \mathbf{Q} 使

$$\text{cond}(\mathbf{PAQ}) < \text{cond}(\mathbf{A}).$$

一般选择 \mathbf{P}, \mathbf{Q} 为对角阵或者三角矩阵.

当矩阵 \mathbf{A} 的元素大小不均时, 对 \mathbf{A} 的行(或列)引进适当的比例因子(使矩阵 \mathbf{A} 的所有行或列按 $\|\cdot\|_1$ 范数大体上有相同的长度,

使 \mathbf{A} 的系数均衡), 对 \mathbf{A} 的条件数是有影响的. 这种方法不能保证 \mathbf{A} 的条件数一定得到改善.

例 10 设

$$\begin{array}{cc} 1 & 10^4 \\ 1 & 1 \end{array} \begin{array}{c} x_1 \\ x_2 \end{array} = \begin{array}{c} 10^4 \\ 2 \end{array}, \quad (6.9)$$

计算 $\text{cond}(\mathbf{A})$.

$$\mathbf{A} = \begin{array}{cc} 1 & 10^4 \\ 1 & 1 \end{array}, \quad \mathbf{A}^{-1} = \frac{1}{10^4 - 1} \begin{array}{cc} 1 & 10^4 \\ 1 & 1 \end{array},$$

$$\text{cond}(\mathbf{A}) = \frac{(1 + 10^4)^2}{10^4 - 1} \approx 10^4.$$

现在 \mathbf{A} 的第一行引进比例因子. 如用 $s_1 = \max_{i=1,2} |a_{i1}| = 10^4$ 除第一个方程式, 得 $\mathbf{A}\mathbf{x} = \mathbf{b}$, 即

$$\begin{array}{cc} 10^{-4} & 1 \\ 1 & 1 \end{array} \begin{array}{c} x_1 \\ x_2 \end{array} = \begin{array}{c} 1 \\ 2 \end{array}, \quad (6.10)$$

而 $(\mathbf{A})^{-1} = \frac{1}{1 - 10^{-4}} \begin{array}{cc} 1 & 1 \\ 1 & 1 - 10^{-4} \end{array},$

于是

$$\text{cond}(\mathbf{A}) = \frac{1}{1 - 10^{-4}} \approx 4.$$

当用列主元消去法解(6.9)时(计算到三位数字),

$$(\mathbf{A} / \mathbf{b}) \quad \left| \begin{array}{cc|c} 1 & 10^4 & 10^4 \\ 0 & -10^4 & -10^4 \end{array} \right.,$$

于是得到很坏的结果: $x_2 = 1$, $x_1 = 0$.

现用列主元消去法解(6.10), 得到

$$(\mathbf{A} / \mathbf{b}) \quad \left| \begin{array}{cc|c} 1 & 1 & 2 \\ 10^{-4} & 1 & 1 \\ 0 & 1 & 1 \end{array} \right.,$$

从而得到较好的计算解: $x_1 = 1$, $x_2 = 1$.

设 $\tilde{\mathbf{x}}$ 为方程组 $\mathbf{A}\mathbf{x} = \mathbf{b}$ 的近似解, 于是可计算 $\tilde{\mathbf{x}}$ 的剩余向量 $\mathbf{r} = \mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}$, 当 \mathbf{r} 很小时, $\tilde{\mathbf{x}}$ 是否为 $\mathbf{A}\mathbf{x} = \mathbf{b}$ 一个较好的近似解呢? 下述定理给出了解答.

定理 24(事后误差估计) 设 \mathbf{A} 为非奇异矩阵, \mathbf{x} 是 $\mathbf{A}\mathbf{x} = \mathbf{b} = \mathbf{0}$ 的精确解. 再设 $\tilde{\mathbf{x}}$ 是此方程组的近似解, $\mathbf{r} = \mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}$, 则

$$\frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|}{\|\mathbf{x}\|} \leq \text{cond}(\mathbf{A}) \cdot \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|}. \quad (6.11)$$

证明 由 $\mathbf{x} - \tilde{\mathbf{x}} = \mathbf{A}^{-1}\mathbf{r}$, 得

$$\mathbf{x} - \tilde{\mathbf{x}} = \mathbf{A}^{-1}\mathbf{r}, \quad (6.12)$$

又有

$$\mathbf{b} = \mathbf{A}\mathbf{x} = \mathbf{A}(\mathbf{x} - \tilde{\mathbf{x}}) + \mathbf{A}\tilde{\mathbf{x}}, \quad \frac{1}{\|\mathbf{x}\|} \|\mathbf{b}\| = \frac{\|\mathbf{A}\|}{\|\mathbf{b}\|}, \quad (6.13)$$

由(6.12)及(6.13)即得到(6.11).

(6.11)式说明, 近似解 $\tilde{\mathbf{x}}$ 的精度(误差界)不仅依赖于剩余 \mathbf{r} 的“大小”, 而且依赖于 \mathbf{A} 的条件数. 当 \mathbf{A} 是病态时, 即使有很小的剩余 \mathbf{r} , 也不能保证 $\tilde{\mathbf{x}}$ 是高精度的近似解.

5.6.2 迭代改善法

设 $\mathbf{A}\mathbf{x} = \mathbf{b}$, 其中 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 为非奇异矩阵, 且为病态方程组(但不过分病态). 当求得方程组的近似解 \mathbf{x} , 下面研究改善方程组近似解 \mathbf{x} 精度的方法.

首先用选主元三角分解法实现分解计算

$$\mathbf{P}\mathbf{A} = \mathbf{L}\mathbf{U},$$

其中 \mathbf{P} 为置换阵, \mathbf{L} 为单位下三角阵, \mathbf{U} 为上三角阵, 且求得计算解 \mathbf{x} .

现利用 \mathbf{x} 的剩余向量来提高 \mathbf{x} 的精度.

计算剩余向量

$$\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k, \quad (6.14)$$

求解 $\mathbf{A}\mathbf{d} = \mathbf{r}_k$, 得到的解记为 \mathbf{d}_k . 然后改善

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{d}_k. \quad (6.15)$$

显然, 如果(6.14), (6.15)及解 $\mathbf{A}\mathbf{d} = \mathbf{r}_k$ 的计算没有误差, 则 \mathbf{x}_{k+1} 就是 $\mathbf{A}\mathbf{x} = \mathbf{b}$ 的精确解. 事实上,

$$\mathbf{A}\mathbf{x}_{k+1} = \mathbf{A}(\mathbf{x}_k + \mathbf{d}_k) = \mathbf{A}\mathbf{x}_k + \mathbf{A}\mathbf{d}_k = \mathbf{A}\mathbf{x}_k + \mathbf{r}_k = \mathbf{b}.$$

但是, 在实际计算中, 由于有舍入误差, \mathbf{x}_{k+1} 只是方程组的近似解, 重复(6.14), (6.15)过程, 就产生一近似解序列 $\{\mathbf{x}_k\}$, 有时可能得到比较好的近似.

算法 5(迭代改善法) 设 $\mathbf{A}\mathbf{x} = \mathbf{b}$, 其中 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 为非奇异矩阵, 且 $\mathbf{A}\mathbf{x} = \mathbf{b}$ 为病态方程组(但不过分病态), 用选主元分解法实现 $\mathbf{PA} = \mathbf{LU}$ 及计算解 \mathbf{x}_k . 本算法用迭代改善法提高近似解 \mathbf{x}_k 精度. 设计算机字长为 t , 用数组 $A(n, n)$ 保存 \mathbf{A} 元素, 数组 $C(n, n)$ 保存三角矩阵 \mathbf{L} 及 \mathbf{U} , 用 $Ip(n)$ 记录行交换信息, $x(n)$ 存贮 \mathbf{x}_k 及 \mathbf{x}_{k+1} , $r(n)$ 保存 \mathbf{r}_k 或 \mathbf{d}_k .

1. 用选主元三角分解实行分解计算

$\mathbf{PA} = \mathbf{LU}$ 且求计算解 \mathbf{x}_k (用单精度)

2. 对于 $k = 1, 2, \dots, N_0$

(1) 计算 $\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k$ (用原始 \mathbf{A} 及双精度计算)

(2) 求解 $\mathbf{L}\mathbf{U}\mathbf{d}_k = \mathbf{P}\mathbf{r}_k$, 即 $\begin{cases} \mathbf{Ly} = \mathbf{Pr}_k \\ \mathbf{Ud}_k = \mathbf{y} \end{cases}$ (用单精度计算)

(3) 如果 $\|\mathbf{d}_k\| / \|\mathbf{x}_k\| > 10^{-t}$ 则输出 $k, \mathbf{x}_k, \mathbf{r}_k$, 停机

(4) 改善 $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{d}_k$ (用单精度计算)

3. 输出迭代改善方法迭代 N_0 次失败信息

当 $\mathbf{A}\mathbf{x} = \mathbf{b}$ 不是过分病态时, 迭代改善法是比较好的改进近似解精度的一种方法, 当 $\mathbf{A}\mathbf{x} = \mathbf{b}$ 非常病态时, $\{\mathbf{x}_k\}$ 可能不收敛.

迭代改善法的实现要依赖于机器及需要保留 \mathbf{A} 的原始副本 .

例 11 用迭代改善法解

$$\begin{array}{rrrr} 1.0303 & 0.99030 & x_1 & = 2.4944 \\ 0.99030 & 0.95285 & x_2 & = 2.3988 \end{array} \quad (\text{记为 } \mathbf{Ax} = \mathbf{b})$$

(这里 $= 10$, $t=5$, 用 5 位浮点数运算) .

解 精确解 $\mathbf{x}^* = (1.2240, 1.2454)^T$ (舍入到小数后第 4 位) .

容易计算

$$\text{cond}(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\| = 2 \times 2000 = 4000.$$

首先实现分解计算 $\mathbf{A} = \mathbf{LU}$, 且求 \mathbf{x} .

$$\mathbf{A} = \begin{matrix} 1 & 0 & 1.0303 & 0.99030 \\ 0.9118 & 1 & 0 & 0.00099 \end{matrix} = \mathbf{LU},$$

且得计算解 $\mathbf{x} = (1.2560, 1.2121)^T$.

应用迭代改善法需要用原始矩阵 \mathbf{A} 且用双倍字长精度计算剩余向量 $\mathbf{r} = \mathbf{b} - \mathbf{Ax}$, 其他计算用单精度 . 计算如下表 .

\mathbf{x}	\mathbf{r}	\mathbf{d}	\mathbf{x}	\mathbf{r}	\mathbf{d}
1.2560	5.7×10^{-7}	- 0.03220	1.2238	1.18×10^{-6}	2.285×10^{-4}
1.2121	3.3715×10^{-5}	0.033502	1.2456	9×10^{-7}	-2.365×10^{-4}

$$\mathbf{x} = (1.2240, 1.2454)^T,$$

$$\mathbf{r} = (-0.682 \times 10^{-5}, -0.659 \times 10^{-5})^T,$$

$$\mathbf{d} = (0.2717 \times 10^{-4}, -0.3515 \times 10^{-4})^T.$$

如果 \mathbf{x} 需要更多的数位, 迭代可以继续 .

5.7 矩阵的正交三角化及应用

本节介绍初等反射阵及平面旋转阵, 矩阵正交约化, 它们在矩

阵计算中起着重要作用 .

5.7.1 初等反射阵

定义 9 设向量 $\mathbf{w} \in \mathbb{R}^n$ 且 $\mathbf{w}^\top \mathbf{w} = 1$, 称矩阵

$$\mathbf{H}(\mathbf{w}) = \mathbf{I} - 2\mathbf{w}\mathbf{w}^\top$$

为初等反射阵(或称为豪斯霍尔德(Householder)变换). 如果记 $\mathbf{w} = (w_1, w_2, \dots, w_n)^\top$, 则

$$\mathbf{H}(\mathbf{w}) = \begin{matrix} 1 - 2w_1^2 & -2w_1 w_2 & \dots & -2w_1 w_n \\ -2w_2 w_1 & 1 - 2w_2^2 & \dots & -2w_2 w_n \\ \dots & \dots & \dots & \dots \\ -2w_n w_1 & -2w_n w_2 & \dots & 1 - 2w_n^2 \end{matrix}.$$

定理 25 设有初等反射阵 $\mathbf{H} = \mathbf{I} - 2\mathbf{w}\mathbf{w}^\top$, 其中 $\mathbf{w}^\top \mathbf{w} = 1$, 则:

(1) \mathbf{H} 是对称矩阵, 即 $\mathbf{H}^\top = \mathbf{H}$.

(2) \mathbf{H} 是正交矩阵, 即 $\mathbf{H}^{-1} = \mathbf{H}$.

(3) 设 \mathbf{A} 为对称矩阵, 那么 $\mathbf{A} = \mathbf{H}^{-1} \mathbf{A} \mathbf{H} = \mathbf{H} \mathbf{A} \mathbf{H}$ 亦是对称矩阵 .

证明 只证 \mathbf{H} 的正交性, 其他显然 .

$$\begin{aligned} \mathbf{H}^\top \mathbf{H} &= \mathbf{H} = (\mathbf{I} - 2\mathbf{w}\mathbf{w}^\top)(\mathbf{I} - 2\mathbf{w}\mathbf{w}^\top) \\ &= \mathbf{I} - 4\mathbf{w}\mathbf{w}^\top + 4\mathbf{w}(\mathbf{w}^\top \mathbf{w})\mathbf{w}^\top = \mathbf{I}. \end{aligned}$$

设向量 $\mathbf{u} \neq \mathbf{0}$, 则显然

$$\mathbf{H} = \mathbf{I} - 2 \frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2}$$

是一个初等反射阵 .

下面考察初等反射阵的几何意义. 考虑以 \mathbf{w} 为法向量且过原点 O 的超平面 S : $\mathbf{w}^\top \mathbf{x} = 0$. 设任意向量 $\mathbf{v} \in \mathbb{R}^n$, 则 $\mathbf{v} = \mathbf{x} + \mathbf{y}$, 其中 $\mathbf{x} \in S$, $\mathbf{y} \perp S$. 于是

$$\mathbf{H}\mathbf{x} = (\mathbf{I} - 2\mathbf{w}\mathbf{w}^\top)\mathbf{x} = \mathbf{x} - 2\mathbf{w}\mathbf{w}^\top \mathbf{x} = \mathbf{x}.$$

对于 $\mathbf{y} \in S$, 易知 $H\mathbf{y} = -\mathbf{y}$, 从而对任意向量 $\mathbf{v} \in \mathbf{R}^n$, 总有

$$H\mathbf{v} = \mathbf{x} - \mathbf{y} = \mathbf{v},$$

其中 \mathbf{v} 为 \mathbf{v} 关于平面 S 的镜面反射(见图 5-1).

初等反射矩阵在计算上的意义

是它能用来约化矩阵, 例如设向量 \mathbf{x}

0. 可选择一初等反射阵 H 使 $H\mathbf{x} = \mathbf{e}_1$. 为此给出下面定理.

定理 26 设 \mathbf{x}, \mathbf{y} 为两个不相等的 n 维向量, $\|\mathbf{x}\|^2 = \|\mathbf{y}\|^2$, 则存在一个初等反射阵 H , 使 $H\mathbf{x} = \mathbf{y}$.

证明 令 $\mathbf{w} = \frac{\mathbf{x} - \mathbf{y}}{\|\mathbf{x} - \mathbf{y}\|}$, 则得到一个初等反射阵

$$H = \mathbf{I} - 2\mathbf{w}\mathbf{w}^T = \mathbf{I} - 2 \frac{(\mathbf{x} - \mathbf{y})}{\|\mathbf{x} - \mathbf{y}\|^2} (\mathbf{x}^T - \mathbf{y}^T),$$

而且

$$\begin{aligned} H\mathbf{x} &= \mathbf{x} - 2 \frac{\mathbf{x} - \mathbf{y}}{\|\mathbf{x} - \mathbf{y}\|^2} (\mathbf{x}^T - \mathbf{y}^T) \mathbf{x} \\ &= \mathbf{x} - 2 \frac{(\mathbf{x} - \mathbf{y})(\mathbf{x}^T \mathbf{x} - \mathbf{y}^T \mathbf{x})}{\|\mathbf{x} - \mathbf{y}\|^2}. \end{aligned}$$

因为

$$\mathbf{x}^T \mathbf{y} = (\mathbf{x} - \mathbf{y})^T (\mathbf{x} - \mathbf{y}) = 2(\mathbf{x}^T \mathbf{x} - \mathbf{y}^T \mathbf{x}),$$

所以

$$H\mathbf{x} = \mathbf{x} - (\mathbf{x} - \mathbf{y}) = \mathbf{y}.$$

容易说明, \mathbf{w} 是使 $H\mathbf{x} = \mathbf{y}$ 成立的唯一长度等于 1 的向量(不计符号).

定理 27(约化定理) 设 $\mathbf{x} = (x_1, x_2, \dots, x_n)^T \neq \mathbf{0}$, 则存在初等反射阵 H 使 $H\mathbf{x} = -\mathbf{e}_1$, 其中

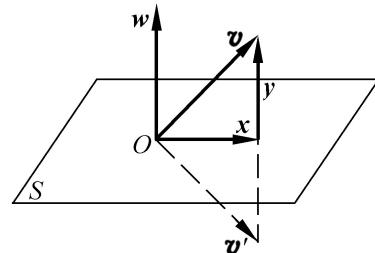


图 5-1

$$\begin{aligned}
 \mathbf{H} &= \mathbf{I} - \frac{1}{x_1^2} \mathbf{u} \mathbf{u}^T, \\
 &= \operatorname{sgn}(x_1) \mathbf{x}_{-2}, \\
 \mathbf{u} &= \mathbf{x} + \mathbf{e}, \\
 &= \frac{1}{2} (\mathbf{u} - \frac{\mathbf{u}}{2}) = (\mathbf{u} + x_1).
 \end{aligned}$$

证明 记 $\mathbf{y} = \mathbf{x} - \mathbf{e}$, 设 $\mathbf{x} = \mathbf{y}$, 取 $= \pm \mathbf{x}_{-2}$, 则有 $\mathbf{x}_{-2} = \mathbf{y}_{-2}$, 于是由定理 22 存在 \mathbf{H} 变换

$$\mathbf{H} = \mathbf{I} - 2\mathbf{w}\mathbf{w}^T,$$

其中 $\mathbf{w} = \frac{\mathbf{x} + \mathbf{e}}{\|\mathbf{x} + \mathbf{e}\|_2}$, 使 $\mathbf{H}\mathbf{x} = \mathbf{y} = \mathbf{x} - \mathbf{e}$.

记 $\mathbf{u} = \mathbf{x} + \mathbf{e} = (u_1, u_2, \dots, u_n)^T$. 于是

$$\mathbf{H} = \mathbf{I} - 2 \frac{\mathbf{u}\mathbf{u}^T}{\|\mathbf{u}\|_2^2} = \mathbf{I} - \frac{1}{u_1^2} \mathbf{u} \mathbf{u}^T,$$

其中 $\mathbf{u} = (x_1 + \dots, x_2, \dots, x_n)^T$, $= \frac{1}{2} (\mathbf{u} - \frac{\mathbf{u}}{2})$.

显然

$$\begin{aligned}
 &= \frac{1}{2} (\mathbf{u} - \frac{\mathbf{u}}{2}) = \frac{1}{2} ((x_1 + \dots)^2 + x_2^2 + \dots + x_n^2) \\
 &= (\mathbf{u} + x_1).
 \end{aligned}$$

如果 x_1 和 x_1 异号, 那么计算 $x_1 +$ 时有效数字可能损失, 我们取 x_1 和 x_1 有相同的符号, 即取

$$= \operatorname{sgn}(x_1) \mathbf{x}_{-2} = \operatorname{sgn}(x_1) \sum_{i=1}^n x_i^2^{1/2}.$$

在计算 $x_1 +$ 时, 可能上溢或下溢, 为了避免溢出, 将 \mathbf{x} 规范化

$$d = \|\mathbf{x}\|, \quad \mathbf{x} = \frac{\mathbf{x}}{d} \quad (\text{设 } d \neq 0),$$

则有 \mathbf{H} 使 $\mathbf{H}\mathbf{x} = \mathbf{e}$, 其中

$$\begin{aligned}
 \mathbf{H} &= \mathbf{I} - (\mathbf{u} + x_1)^{-1} \mathbf{u} \mathbf{u}^T, \\
 &= / d, \quad \mathbf{u} = \mathbf{u} / d, \quad = / d^2, \\
 \mathbf{H} &= \mathbf{H}.
 \end{aligned}$$

算法 6 设 $\mathbf{x} = (x_1, x_2, \dots, x_n)^T \neq \mathbf{0}$, 本算法计算 \mathbf{u} , 及 使 $(\mathbf{I} - (\mathbf{u}\mathbf{u}^T)^{-1})\mathbf{x} = \mathbf{e}$, \mathbf{u} 的分量冲掉 \mathbf{x} 的分量.

$$(1) \text{ 计算 } d \quad d = \max_{1 \leq i \leq n} |x_i|$$

$$(2) \quad x_i \quad u_i = x_i / d \quad (i = 1, 2, \dots, n)$$

$$(3) \quad \text{sgn}(u_i) = \begin{cases} + & u_i \geq 0 \\ - & u_i < 0 \end{cases}$$

$$(4) \quad u_i = u_i +$$

$$(5) \quad u_i^* = u_i$$

$$(6) \quad d^* =$$

关于 \mathbf{HA} 的计算, 设 $\mathbf{A} = (\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_n) \in \mathbf{R}^{m \times n}$, 其中 \mathbf{a}_i 为 \mathbf{A} 的第 i 列向量, 则

$$\mathbf{HA} = (\mathbf{H}\mathbf{a}_1 \ \mathbf{H}\mathbf{a}_2 \ \dots \ \mathbf{H}\mathbf{a}_n),$$

因此计算 \mathbf{HA} 就是要计算

$$\mathbf{H}\mathbf{a}_i = (\mathbf{I} - (\mathbf{u}\mathbf{u}^T)^{-1})\mathbf{a}_i = \mathbf{a}_i - ((\mathbf{u}\mathbf{u}^T)^{-1}\mathbf{u}\mathbf{a}_i)\mathbf{u} \quad (i = 1, 2, \dots, n).$$

于是计算 $\mathbf{H}\mathbf{a}_i$ 只需计算两向量的数量积和两向量的加法. 计算 \mathbf{HA} 只需作 $2nm$ 次乘法运算.

5.7.2 平面旋转矩阵

设 $\mathbf{x}, \mathbf{y} \in \mathbf{R}^2$, 则变换

$$\begin{matrix} y_1 \\ y_2 \end{matrix} = \begin{matrix} \cos & \sin \\ -\sin & \cos \end{matrix} \begin{matrix} x_1 \\ x_2 \end{matrix}, \quad \text{或} \quad \mathbf{y} = \mathbf{Px}$$

是平面上向量的一个旋转变换, 其中

$$\mathbf{P}(\theta) = \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix}$$

为正交矩阵.

\mathbf{R}^n 中变换: $\mathbf{y} = \mathbf{Px}$,

其中 $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$, $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$, 而

$$\begin{array}{ccccccc}
 & & i & & j & & \\
 & & 1 & & & & \\
 & & w & & & & \\
 & & 1 & & & & \\
 & & \cos & & \sin & & i \\
 & & & & & & \\
 & & & & & 1 & \\
 \mathbf{P} = \mathbf{P}(i, j, \quad) = & & & & & w & \\
 & & & & & -\sin & \\
 & & & & & \cos & j \\
 & & & & & & \\
 & & & & & 1 & \\
 & & & & & w & \\
 & & & & & 1 &
 \end{array}$$

称为 \mathbf{R}^n 中平面 $\{x_i, x_j\}$ 的旋转变换(或称为吉文斯(Givens)变换), $\mathbf{P} = \mathbf{P}(i, j, \quad) = P(i, j)$ 称为平面旋转矩阵.

显然, $\mathbf{P}(i, j, \quad)$ 具有性质:

(1) \mathbf{P} 与单位阵 \mathbf{I} 只是在 $(i, i), (i, j), (j, i), (j, j)$ 位置元素不一样, 其他相同.

(2) \mathbf{P} 为正交矩阵 ($\mathbf{P}^{-1} = \mathbf{P}^T$).

(3) $\mathbf{P}(i, j) \mathbf{A}$ (左乘) 只需计算第 i 行与第 j 行元素, 即对 $\mathbf{A} = (a_{ij})_{m \times n}$ 有

$$\begin{array}{ccc}
 a_{il} & = & c & s & a_{il} \\
 a_{jl} & = & -s & c & a_{jl}
 \end{array} \quad (l = 1, 2, \dots, n).$$

其中 $c = \cos, s = \sin$.

(4) $\mathbf{AP}(i, j)$ (右乘) 只需计算第 i 列与第 j 列元素

$$(a_{i1}, a_{ij}) = (a_{i1}, a_{lj}) \begin{pmatrix} c & s \\ -s & c \end{pmatrix} \quad (l = 1, 2, \dots, m).$$

利用平面旋转变换, 可使向量 \mathbf{x} 中的指定元素变为零.

定理28(约化定理) 设 $\mathbf{x} = (x_1, \dots, x_i, \dots, x_j, \dots, x_n)^T$, 其中 x_i, x_j 不全为零, 则可选择平面旋转阵 $\mathbf{P}(i, j,)$, 使

$$i \quad j$$

$$\mathbf{P}\mathbf{x} = (x_1, \dots, x_i, \dots, 0, \dots, x_n)^T,$$

其中 $x_i = \sqrt{x_i^2 + x_j^2}$, $= \arctan(x_j/x_i)$.

证明 取 $c = \cos = x_i/x_i$, $s = \sin = x_j/x_i$. 由 $\mathbf{P}(i, j,)\mathbf{x} = \mathbf{x} = (x_1, x_2, \dots, x_i, \dots, x_j, \dots, x_n)^T$, 利用矩阵乘法, 显然有

$$x_i = cx_i + sx_j,$$

$$x_j = -sx_i + cx_j,$$

$$x_k = x_k \quad (k \neq i, j).$$

于是, 由 c, s 的取法得

$$x_i = \sqrt{x_i^2 + x_j^2}, \quad x_j = 0.$$

5.7.3 矩阵的 QR 分解

下面讨论用正交矩阵来约化矩阵, 可得到下述结果.

设 $\mathbf{A} \in \mathbf{R}^{m \times n}$ 且为非零矩阵, 则存在初等反射阵 $\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_s$ 使

$$\mathbf{H}_s \dots \mathbf{H}_2 \mathbf{H}_1 \mathbf{A} = \mathbf{A}^{(s+1)} \quad (\text{上梯形}).$$

设有

$$\begin{aligned} & a_{11} \quad a_{12} \quad \dots \quad a_{1n} \\ \mathbf{A} = \mathbf{A}^{(1)} = & \begin{array}{cccc} a_{11} & a_{12} & \dots & a_{1n} \\ \cdots & \cdots & \cdots & \cdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{array} = (\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_n) \quad (\text{按列分块}). \end{aligned} \tag{7.1}$$

(1) 第1步约化: 如果 $\mathbf{a}_1 = \mathbf{0}$, 取 $\mathbf{H}_1 = \mathbf{I}$, 即这一步不需要约化, 不妨设 $\mathbf{a}_1 \neq \mathbf{0}$, 于是可选取初等反射阵 $\mathbf{H}_1 = \mathbf{I} - \frac{1}{\| \mathbf{a}_1 \|} \mathbf{u}_1 \mathbf{u}_1^T$ 使

$$\mathbf{H} \mathbf{a} = -\mathbf{e}_1.$$

于是

$$\begin{aligned} \mathbf{H} \mathbf{A}^{(1)} &= (\mathbf{H} \mathbf{a} \quad \mathbf{H} \mathbf{a} \quad \dots \quad \mathbf{H} \mathbf{a}) \\ &= \begin{matrix} -1 & a_{12}^{(2)} & \dots & a_{1n}^{(2)} \\ 0 & a_{22}^{(2)} & \dots & a_{2n}^{(2)} \\ \dots & \dots & \dots & \dots \\ 0 & a_{m2}^{(2)} & \dots & a_{mn}^{(2)} \end{matrix} = \begin{matrix} -1 & r_2 & \mathbf{B} \\ \mathbf{0} & \mathbf{c} & \mathbf{D} \end{matrix}, \end{aligned}$$

其中 $\mathbf{c} = (a_{22}^{(2)}, \dots, a_{m2}^{(2)})^T \in \mathbf{R}^{m-1}$, $\mathbf{D} \in \mathbf{R}^{(m-1) \times (n-2)}$.

(2) 第 k 步约化: 设已完成对 \mathbf{A} 上述第 1 步...第 $k-1$ 步的约化, 再进行第 k 步约化. 即存在初等反射阵 $\mathbf{H}_k, \mathbf{H}_{k-1}, \dots, \mathbf{H}_1$ 使

$$\mathbf{H}_{k-1} \dots \mathbf{H}_1 \mathbf{H} \mathbf{A} = \mathbf{A}^{(k)},$$

其中

$$\begin{aligned} \mathbf{A}^{(k)} &= \begin{matrix} -1 & a_{12}^{(2)} & \dots & a_{1k-1}^{(2)} & a_{1k}^{(2)} & \dots & a_{1n}^{(2)} \\ -2 & \dots & a_{2k-1}^{(3)} & a_{2k}^{(3)} & \dots & a_{2n}^{(3)} \\ \vdots & \dots & \dots & \dots & \dots & \dots & \dots \\ k-1 & a_{kk}^{(k)} & \dots & a_{k-1k}^{(k-1)} & a_{k-1k}^{(k-1)} & \dots & a_{k-1n}^{(k-1)} \\ & & & & & & a_{kn}^{(k)} \\ & & & & & & \dots \\ & & & & a_{mk}^{(k)} & \dots & a_{mn}^{(k)} \end{matrix} \\ &= \begin{matrix} \mathbf{R}_k & \mathbf{r}_k & \mathbf{B}_k & k-1 \\ \mathbf{0} & \mathbf{c}_k & \mathbf{D}_k & m-k+1 \end{matrix}, \end{aligned}$$

其中, \mathbf{R}_k 为 $k-1$ 阶上三角阵, $\mathbf{c}_k \in \mathbf{R}^{m-k+1}$, $\mathbf{D}_k \in \mathbf{R}^{(m-k+1) \times (n-k)}$.

不妨设 $\mathbf{c}_k = \mathbf{0}$, 否则这一步不需要约化(如果 \mathbf{A} 列满秩, 则 $\mathbf{c}_k \neq \mathbf{0}$). 于是, 可选取初等反射阵 $\mathbf{H}_k = \mathbf{I}_{m-k+1} - \frac{1}{\|u_k\|^2} u_k u_k^T$ 使

$$\mathbf{H}_k \mathbf{c}_k = -\mathbf{e}_k.$$

令

$$\begin{aligned} \mathbf{H}_k &= \begin{array}{cc} \mathbf{I}_{k-1} & k-1 \\ & \mathbf{H}_k \quad m-k+1 \end{array}, \end{aligned}$$

第 k 步约化为

$$\begin{aligned} \mathbf{H}_k \mathbf{A}^{(k)} &= \mathbf{H}_k \mathbf{H}_{k-1} \dots \mathbf{H}_1 \mathbf{H}_1 \mathbf{A} = \mathbf{A}^{(k+1)} \\ &= \begin{array}{ccccc} \mathbf{I}_{k-1} & \mathbf{0} & \mathbf{R}_k & \mathbf{r}_k & \mathbf{B}_k \\ = & \mathbf{0} & \mathbf{H}_k & \mathbf{0} & \mathbf{c}_k & \mathbf{D}_k \\ & \mathbf{R}_k & \mathbf{r}_k & \mathbf{B}_k \\ = & \mathbf{0} & \mathbf{H}_k \mathbf{c}_k & \mathbf{H}_k \mathbf{D}_k \end{array}, \end{aligned}$$

其中 $\mathbf{A}^{(k+1)}$ 左上角 k 阶子矩阵为 k 阶上三角阵 \mathbf{R}_{k+1} , 这就使 \mathbf{A} 三
角化过程前进了一步.

令 $s = \min(m-1, n)$, 继续上述过程, 最后有

$$\mathbf{H}_s \dots \mathbf{H}_1 \mathbf{H}_1 \mathbf{A} = \mathbf{A}^{(s+1)} = \mathbf{R}(\text{上梯形}).$$

第 k 步需要计算 $(_{k-1}, \mathbf{u}_k, {}_{k-1})$ 及 $\mathbf{H}_k \mathbf{D}_k$, 第 k 步约化大约需要
 $2 \cdot (m-k+1)(n-k)$ 次乘法运算.

总结上述讨论给出下述结果.

定理 29(矩阵的正交约化) 设 $\mathbf{A} \in \mathbb{R}^{m \times n}$ 且 $\mathbf{A} \neq \mathbf{0}$, $s = \min(m-1, n)$, 则存在初等反射阵 $\mathbf{H}_s, \mathbf{H}_{s-1}, \dots, \mathbf{H}_1$ 使

$$\mathbf{H}_s \dots \mathbf{H}_1 \mathbf{H}_1 \mathbf{A} = \mathbf{R}(\text{上梯形}),$$

且计算量约为 $n^2 m - n^3 / 3$ (当 $m \geq n$) 次乘法运算.

定理 30(矩阵的 QR 分解)

(1) 设 $\mathbf{A} \in \mathbb{R}^{m \times n}$ 且 \mathbf{A} 的秩为 n ($m > n$), 则存在初等反射阵 $\mathbf{H}_s, \mathbf{H}_{s-1}, \dots, \mathbf{H}_1$ 使

$$\mathbf{H}_s \dots \mathbf{H}_1 \mathbf{H}_1 \mathbf{A} = \begin{array}{c} \mathbf{R} \\ \mathbf{0} \end{array},$$

其中 \mathbf{R} 为 n 阶非奇异上三角阵.

(2) 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 为非奇异矩阵, 则 \mathbf{A} 有分解

$$\mathbf{A} = \mathbf{Q}\mathbf{R},$$

其中 \mathbf{Q} 为正交矩阵, \mathbf{R} 为上三角矩阵. 且当 \mathbf{R} 具有正对角元时, 分解唯一.

证明 (1) 由定理 29 可得.

(2) 由设及定理 29 存在初等反射阵 $\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_{n-1}$ 使

$$\begin{array}{ccccccccc} & r_{11} & r_{12} & \cdots & r_{1n} & & & & \\ & & r_{22} & \cdots & r_{2n} & & & & \\ \mathbf{H}_{n-1} \dots \mathbf{H}_2 \mathbf{H}_1 \mathbf{A} = & & & & & & & & \mathbf{R}. \\ & & & & & & \ddots & & \\ & & & & & & & \ddots & \\ & & & & & & & & r_{nn} \end{array}$$

记 $\mathbf{Q}^T = \mathbf{H}_{n-1} \dots \mathbf{H}_2 \mathbf{H}_1$, 则上式为

$$\mathbf{Q}^T \mathbf{A} = \mathbf{R},$$

即

$$\mathbf{A} = \mathbf{Q} \mathbf{R},$$

其中 \mathbf{Q} 为正交矩阵, \mathbf{R} 为上三角阵.

唯一性. 设有 $\mathbf{A} = \mathbf{Q}_1 \mathbf{R}_1 = \mathbf{Q}_2 \mathbf{R}_2$, 其中 $\mathbf{Q}_1, \mathbf{Q}_2$ 为正交阵, $\mathbf{R}_1, \mathbf{R}_2$ 为非奇异上三角阵, 且 $\mathbf{R}_1, \mathbf{R}_2$ 具有正对角元素, 则

$$\mathbf{A}^T \mathbf{A} = \mathbf{R}_1^T \mathbf{Q}_1^T \mathbf{Q}_1 \mathbf{R}_1 = \mathbf{R}_1^T \mathbf{R}_1,$$

$$\mathbf{A}^T \mathbf{A} = \mathbf{R}_2^T \mathbf{Q}_2^T \mathbf{Q}_2 \mathbf{R}_2 = \mathbf{R}_2^T \mathbf{R}_2.$$

由假设及对称正定矩阵 $\mathbf{A}^T \mathbf{A}$ 的 Cholesky 分解的唯一性, 则 $\mathbf{R}_1 = \mathbf{R}_2$. 从而可得 $\mathbf{Q}_1 = \mathbf{Q}_2$.

下面考虑用平面旋转变换来约化矩阵.

定理 31 (用吉文斯变换计算矩阵的 QR 分解) 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 为非奇异矩阵, 则

(1) 存在正交矩阵 $\mathbf{P}_1, \mathbf{P}_2, \dots, \mathbf{P}_{n-1}$ 使

$$\begin{array}{ccccccccc} & r_{11} & r_{12} & \cdots & r_{1n} & & & & \\ & & r_{22} & \cdots & r_{2n} & & & & \\ \mathbf{P}_{n-1} \dots \mathbf{P}_2 \mathbf{P}_1 \mathbf{A} = & & & & & & & & \mathbf{R}. \\ & & & & & & \ddots & & \\ & & & & & & & \ddots & \\ & & & & & & & & r_{nn} \end{array}$$

(2) \mathbf{A} 有 QR 分解

$$\mathbf{A} = \mathbf{QR},$$

其中 \mathbf{Q} 为正交阵, \mathbf{R} 为非奇异上三角阵, 且当 \mathbf{R} 对角元素都为正时, 分解是唯一的.

证明 (1) 第 1 步约化: 由设有 $j(1 \quad j \quad n)$ 使 $a_{j1} = 0$. 若 $a_{j1} \neq 0$ ($j = 2, \dots, n$), 则可选择吉文斯变换 $\mathbf{P}(1, j)$ ($j = 2, \dots, n$) 使

$$\begin{array}{ccccccccc} & r_{11} & r_{12} & \cdots & r_{1n} & & & & \\ & & a_{22}^{(2)} & \cdots & a_{nn}^{(2)} & & & & \\ \mathbf{P}(1, n) \dots \mathbf{P}(1, 2) \mathbf{A} = & & & & & & & & \mathbf{A}^{(2)}, \\ & & & & & & & \cdots & \\ & & a_{n2}^{(2)} & \cdots & a_{nn}^{(2)} & & & & \end{array}$$

可简记为 $\mathbf{P}_1 \mathbf{A} = \mathbf{A}^{(2)}$, 其中 $\mathbf{P} = \mathbf{P}(1, n) \dots \mathbf{P}(1, 2)$.

(2) 第 k 步约化: 设上述过程已完成第 1 步至第 $k-1$ 步, 于是有

$$\begin{array}{ccccccccc} & r_{11} & r_{12} & \cdots & r_{1k} & \cdots & r_{1n} & & \\ & r_{21} & \cdots & r_{2k} & \cdots & r_{2n} & & & \\ & & & & & & & & \\ & & & & w & \cdots & \cdots & & \\ \mathbf{P}_{k-1} \dots \mathbf{P}_2 \mathbf{P}_1 \mathbf{A} = & & & & a_{kk}^{(k)} & \cdots & a_{kn}^{(k)} & & \mathbf{A}^{(k)} \\ & & & & & & & & \\ & & & & & & & & \\ & & & & a_{nk}^{(k)} & \cdots & a_{nn}^{(k)} & & \end{array}.$$

由设有 $j(n \quad j \quad k)$ 使 $a_{jk}^{(k)} = 0$, 若 $a_{jk}^{(k)} \neq 0$ ($j = k+1, \dots, n$), 则可选择吉文斯变换 $\mathbf{P}(k, j)$ ($j = k+1, \dots, n$) 使

$$\begin{aligned} \mathbf{P}_k \mathbf{A}^{(k)} &= \mathbf{P}(k, n) \dots \mathbf{P}(k, k+1) \mathbf{A}^{(k)} \\ &= \mathbf{P}_k \mathbf{P}_{k-1} \dots \mathbf{P}_1 \mathbf{A} = \mathbf{A}^{(k+1)}, \end{aligned}$$

其中 $\mathbf{P}_k = \mathbf{P}(k, n) \dots \mathbf{P}(k, k+1)$.

(3) 继续上述约化过程, 最后则有

$$\mathbf{P}_{n-1} \dots \mathbf{P} \mathbf{P} \mathbf{A} = \mathbf{R} \quad (\text{上三角阵}),$$

其中 \mathbf{P}_k 为正交阵(为一系列平面旋转阵的乘积). 记 $\mathbf{Q}^T = \mathbf{P}_{n-1} \dots$

\mathbf{P} , 则 \mathbf{A} 有 QR 分解

$$\mathbf{A} = \mathbf{Q}\mathbf{R},$$

其中 \mathbf{Q} 为正交矩阵, \mathbf{R} 为非奇异上三角阵.

用 \mathbf{H} 变换实现 \mathbf{A} 的正交三角约化需要 $\frac{2}{3}n^3$ 次乘法, n 次开方运算, 而用吉文斯变换约化计算约需要 $\frac{4}{3}n^3$ 次乘法, $\frac{n^2}{2}$ 次开方运算. 这说明, 对一般矩阵 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 用吉文斯变换实现 \mathbf{A} 正交三角约化比用 \mathbf{H} 变换实现 \mathbf{A} 正交三角约化, 计算量要大一倍. 但是, 如果 \mathbf{A} 为三对角阵或上豪斯霍尔德阵, 利用吉文斯变换实现 \mathbf{A} 正交三角约化要简单、经济.

5.7.4 求解超定方程组

设 $\mathbf{Ax} = \mathbf{b}$, 其中 $\mathbf{A} \in \mathbb{R}^{m \times n}$. 当 $m > n$ 时称为超定方程组, 一般地, 没有通常意义上的解.

线性最小二乘问题: 对超定方程组, 寻求 $\mathbf{x}^* \in \mathbb{R}^n$ 使

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{b} - \mathbf{Ax}\|_2^2 = \|\mathbf{b} - \mathbf{Ax}^*\|_2^2.$$

如果使上式成立的 $\mathbf{x}^* \in \mathbb{R}^n$ 存在, 称 \mathbf{x}^* 为超定方程组最小二乘解.

记残差向量 $\mathbf{r} = \mathbf{b} - \mathbf{Ax}$, 于是

$$\|\mathbf{r}\|_2^2 = \|\mathbf{b} - \mathbf{Ax}\|_2^2 = \sum_{i=1}^m b_i^2 - \sum_{j=1}^n a_{ij} x_j^2,$$

即寻求 $\mathbf{x}^* = (x_1^*, \dots, x_n^*)^T \in \mathbb{R}^n$ 使偏差 r_i 的平方和最小.

设 $\mathbf{A} \in \mathbb{R}^{m \times n}$ ($m > n$) 且 \mathbf{A} 具有线性无关的列, 可利用 \mathbf{A} 的正交约化来求超定方程组的最小二乘解.

由正交约化定理, 可选择初等反射阵 $\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_n$ 使

$$\mathbf{H}_1 \cdots \mathbf{H}_n \mathbf{A} = \begin{pmatrix} \mathbf{R} & \\ \mathbf{0} & m-n \end{pmatrix},$$

且同时约化常数项

$$\mathbf{H}_n \dots \mathbf{H}_1 \mathbf{H} \mathbf{b} = \frac{\mathbf{c}}{\mathbf{d}}^T, \quad \mathbf{d} \in \mathbb{R}^{m-n},$$

其中, $\mathbf{R} \in \mathbb{R}^{n \times n}$ 为非奇异矩阵, $\mathbf{c} \in \mathbb{R}^n$.

令 $\mathbf{Q} = \mathbf{H}_n \dots \mathbf{H}_1 \mathbf{H}$, 于是

$$\mathbf{Q}\mathbf{r} = \mathbf{Q}\mathbf{b} - \mathbf{Q}\mathbf{A}\mathbf{x} = \frac{\mathbf{c}}{\mathbf{d}}^T - \frac{\mathbf{R}}{\mathbf{0}} \mathbf{x} = \frac{\mathbf{c} - \mathbf{R}\mathbf{x}}{\mathbf{d}}.$$

因为 \mathbf{Q} 为正交矩阵, 所以

$$\|\mathbf{r}\|_2^2 = \|\mathbf{Q}\mathbf{r}\|_2^2 = \|\mathbf{c} - \mathbf{R}\mathbf{x}\|_2^2 + \|\mathbf{d}\|_2^2 \quad (\text{对任何 } \mathbf{x} \in \mathbb{R}^n).$$

从而, 当选取 \mathbf{x}^* 为 $\mathbf{R}\mathbf{x} = \mathbf{c}$ 的解时, 则

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{r}\|_2^2 = \|\mathbf{b} - \mathbf{A}\mathbf{x}^*\|_2^2 = \|\mathbf{r}^*\|_2^2,$$

且达到极小的残差向量:

$$\begin{aligned} \mathbf{Q}\mathbf{r}^* &= \frac{\mathbf{c} - \mathbf{R}\mathbf{x}^*}{\mathbf{d}} = \frac{\mathbf{0}}{\mathbf{d}}, \\ \mathbf{r}^* &= \mathbf{H}_1 \mathbf{H}_2 \dots \mathbf{H}_n \frac{\mathbf{0}}{\mathbf{d}}. \end{aligned}$$

算法 7(超定方程组) 设 $\mathbf{A}\mathbf{x} = \mathbf{b}$, 其中 $\mathbf{A} \in \mathbb{R}^{m \times n}$ ($m > n$), 且设 $\text{rank}(\mathbf{A}) = n$ (列满秩), 本算法利用正交三角约化 $\mathbf{H}_n \dots \mathbf{H}_1 \mathbf{H} \mathbf{A} = \mathbf{R}$, 其中 $\mathbf{R} \in \mathbb{R}^{n \times n}$ 为非奇异阵, 求 $\mathbf{x}^* \in \mathbb{R}^n$ 使 $\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2 = \|\mathbf{b} - \mathbf{A}\mathbf{x}^*\|_2^2$. 达到极小的残差向量冲掉 \mathbf{b} .

(1) 正交约化

对于 $k = 1, 2, \dots, n$

1) $\mathbf{A} \leftarrow \mathbf{H}_k \mathbf{A}$ (即计算 $\mathbf{D}_k \leftarrow \mathbf{H}_k \mathbf{D}_k$)

2) $\mathbf{b} \leftarrow \mathbf{H}_k \mathbf{b}$

(2) $c_i \leftarrow b_i$ ($i = 1, 2, \dots, n$)

(3) 求解三角形方程组, 求最小二乘解 \mathbf{x}^*

$$\mathbf{R}\mathbf{x} = \mathbf{c}$$

(4) $b_i = 0$ ($i = 1, 2, \dots, n$)

(5) 求达到极小的残差向量

对于 $k = n, n-1, \dots, 1$

$$\mathbf{b} - \mathbf{H}_k \mathbf{b}$$

算法 7 是解线性最小二乘问题的一个非常稳定的方法.

例 12 用正交约化方法求超定方程组

$$\begin{array}{rcl} 2 & -1 & 1 \\ 8 & 4 & 0 \\ 2 & 1 & x_1 = 1 \\ 7 & -1 & x_2 = 8 \\ 4 & 0 & 3 \end{array}$$

的最小二乘解.

解 记 $(\mathbf{A} | \mathbf{b}) = (\mathbf{a} \ \mathbf{a} \ \mathbf{a})$, $\mathbf{A} \in \mathbb{R}^{5 \times 2}$, $\text{rank}(\mathbf{A}) = 2$.

(1) 正交约化

(a) 利用算法 6, 选择 $\mathbf{H} = \mathbf{I} - \mathbf{u} \mathbf{u}^T$, 使 $\mathbf{H} \mathbf{a} = \mathbf{e}_1$,

其中

$$x_1 = 1.463087,$$

$$\mathbf{u} = (1.713087, 1, 0.25, 0.875, 0.5)^T,$$

$$x_2 = x_1 (0.25 + 1) = 2.506395,$$

$$x_3 = 11.704696.$$

于是

$$-11.704696 \quad -2.135895 \quad -6.151378$$

$$3.336931 \quad -4.174556$$

$$\mathbf{H}(\mathbf{A} | \mathbf{b}) = \begin{matrix} 0.8342327 & -0.043639 \\ -1.580185 & 4.347264 \\ -0.331535 & 0.912722 \end{matrix}.$$

$$-1.580185 \quad 4.347264$$

$$-0.331535 \quad 0.912722$$

(b) 对于 $\mathbf{c} = (3.336931, 0.8342327, -1.580185, -0.331535)^T$ 利用算法 6, 确定 $\mathbf{H} = \mathbf{I} - \mathbf{u}_2 \mathbf{u}_2^T$ 使 $\mathbf{H} \mathbf{c} = -\mathbf{u}_2 \mathbf{e}$.

$$\mathbf{u}_2 = 1.138689,$$

$$\mathbf{u}_2 = (2.138689, 0.249999, -0.4735444, -0.09935326)^T,$$

$$\mathbf{u}_2 = 2.435302,$$

$$\mathbf{u}_2 = 3.799727.$$

于是

$$\begin{aligned} & -11.704696 \quad -2.135895 \quad -6.151378 \\ & \quad 0 \quad -3.799727 \quad 5.563212 \\ \mathbf{H} \mathbf{H} (\mathbf{A} / \mathbf{b}) = & \quad 1.094643 \\ & \quad 2.191146 \\ & \quad 0.460352 \\ & = \frac{\mathbf{R} \quad \mathbf{c}}{\mathbf{0} \quad \mathbf{d}}. \end{aligned}$$

(2) 求解 $\mathbf{Rx} = \mathbf{c}$ 即

$$\begin{array}{cccccc} -11.704696 & -2.135895 & x_1 & = & -6.151378 \\ 0 & -3.799727 & x_2 & = & 5.563212 \end{array},$$

得到超定方程组的最小二乘解

$$\mathbf{x}^* = \begin{array}{c} 0.792721 \\ -1.46108 \end{array}.$$

评注

对于良态问题, 高斯消去法也可能给出很坏的结果, 即说明这个算法是不稳定的. 在高斯消去法中引进选主元的技巧, 就得到了解方程组的完全主元素消去法和列主元素消去法, 选主元素技

巧的根本作用是为了对增长因子(即 $r = \max_{1 \leq k \leq n} |a_k| / a$, 其中 $a_k = \max_{1 \leq i, j \leq n} |d_{ij}^{(k)}|$, $a = \max_{1 \leq i, j \leq n} |a_{ij}|$, $\mathbf{A}^{(k)} = (d_{ij}^{(k)})$)进行控制, 即对舍入误差的增长加以控制。由此, 完全选主素消去法及列主元消去法是数值稳定的算法。这两种方法都是计算机上解线性方程组的有效方法。但通常用列主元消去法即可。

从代数上看, 直接分解法和高斯消去法本质上一样, 但如果采用“双精度累加”计算 $a b_i$, 那么直接三角分解法的精度要比高斯消去法为高。

对于对称正定方程组, 采用不选主元素的平方根法(或改进的平方根法)求解适宜。理论分析指出, 解对称正定方程组的平方根法是一个稳定的算法, 在工程计算中使用比较广泛。

追赶法是解三对角线方程组(对角元占优势)的有效方法, 它具有计算量少, 方法简单, 算法稳定等优点。

关于矩阵的条件数, 病态方程组, 算法的稳定性等都是计算数学中比较重要的概念, 在这里我们只做了简单的介绍。本章的内容可进一步参考文献[3], [7], [13]。

习 题

- 设 \mathbf{A} 是对称阵且 $a_{11} \neq 0$, 经过高斯消去法一步后, \mathbf{A} 约化为

$$a_{11} \quad \mathbf{a}^T$$

$$\mathbf{0} \quad \mathbf{A}_2$$

证明 \mathbf{A}_2 是对称矩阵。

- 设 $\mathbf{A} = (a_{ij})_n$ 是对称正定矩阵, 经过高斯消去法一步后, \mathbf{A} 约化为

$$a_{11} \quad \mathbf{a}^T$$

$$\mathbf{0} \quad \mathbf{A}_2,$$

其中 $\mathbf{A}_2 = (d_{ij}^{(2)})_{n-1}$ 。证明:

(1) \mathbf{A} 的对角元素 $a_{ii} > 0$ ($i = 1, 2, \dots, n$) ;

(2) \mathbf{A} 是对称正定矩阵 .

3. 设 \mathbf{L}_k 为指标为 k 的初等下三角阵(除第 k 列对角元以下元素外, \mathbf{L}_k 和单位阵 \mathbf{I} 相同), 即

$$\mathbf{L}_k = \begin{matrix} & & k \\ & 1 & \\ w & & \\ & 1 & \\ m_{k+1,k} & & 1 \\ & \cdots & w \\ & m_{n,k} & 1 \end{matrix}.$$

求证当 $i, j > k$ 时, $\mathbf{L}_k = \mathbf{I}_{ij} \mathbf{L}_k \mathbf{I}_{ij}$ 也是一个指标为 k 的初等下三角阵, 其中 \mathbf{I}_{ij} 为初等置换阵 .

4. 试推导矩阵 \mathbf{A} 的 Crout 分解 $\mathbf{A} = \mathbf{LU}$ 的计算公式, 其中 \mathbf{L} 为下三角阵, \mathbf{U} 为单位上三角阵 .

5. 设 $\mathbf{Ux} = \mathbf{d}$, 其中 \mathbf{U} 为三角矩阵 .

(a) 就 \mathbf{U} 为上及下三角矩阵推导一般的求解公式, 并写出算法 .

(b) 计算解三角形方程组 $\mathbf{Ux} = \mathbf{d}$ 的乘除法次数 .

(c) 设 \mathbf{U} 为非奇异阵, 试推导求 \mathbf{U}^{-1} 的计算公式 .

6. 证明: (a) 如果 \mathbf{A} 是对称正定阵, 则 \mathbf{A}^{-1} 也是对称正定阵;

(b) 如果 \mathbf{A} 是对称正定阵, 则 \mathbf{A} 可唯一地写成 $\mathbf{A} = \mathbf{L}^T \mathbf{L}$, 其中 \mathbf{L} 是具有正对角元的下三角阵 .

7. 用高斯 - 若当方法求 \mathbf{A} 的逆矩阵, 其中

$$\mathbf{A} = \begin{matrix} 2 & 1 & -3 & -1 \\ 3 & 1 & 0 & 7 \\ -1 & 2 & 4 & -2 \\ 1 & 0 & -1 & 5 \end{matrix}.$$

8. 用追赶法解三对角方程组 $\mathbf{Ax} = \mathbf{b}$, 其中

$$\mathbf{A} = \begin{matrix} 2 & -1 & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 \\ 0 & -1 & 2 & -1 & 0 \\ 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & -1 & 2 \end{matrix}, \quad \mathbf{b} = \begin{matrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix}.$$

9. 用改进的平方根法解方程组

$$\begin{matrix} 2 & -1 & 1 & x_1 & 4 \\ -1 & -2 & 3 & x_2 & 5 \\ 1 & 3 & 1 & x_3 & 6 \end{matrix}.$$

10. 下述矩阵能否分解为 \mathbf{LU} (其中 \mathbf{L} 为单位下三角阵, \mathbf{U} 为上三角阵)?

若能分解, 那么分解是否唯一?

$$\mathbf{A} = \begin{matrix} 1 & 2 & 3 \\ 2 & 4 & 1 \\ 4 & 6 & 7 \end{matrix}, \quad \mathbf{B} = \begin{matrix} 1 & 1 & 1 \\ 2 & 2 & 1 \\ 3 & 3 & 1 \end{matrix}, \quad \mathbf{C} = \begin{matrix} 1 & 2 & 6 \\ 2 & 5 & 15 \\ 6 & 15 & 46 \end{matrix}.$$

11. 设

$$\mathbf{A} = \begin{matrix} 0.6 & 0.5 \\ 0.1 & 0.3 \end{matrix},$$

计算 \mathbf{A} 的行范数, 列范数, 2 -范数及 F -范数.

12. 求证: (a) $\|\mathbf{x}\| = \|\mathbf{x}_1\| + \|\mathbf{x}_2\| + \|\mathbf{x}_3\|$,

$$(b) \frac{1}{n} \|\mathbf{A}\|_F = \|\mathbf{A}\|_2 = \|\mathbf{A}\|_F.$$

13. 设 $\mathbf{P} \in \mathbf{R}^{n \times n}$ 且非奇异, 又设 $\|\mathbf{x}\|$ 为 \mathbf{R}^n 上一向量范数, 定义

$$\|\mathbf{x}\|_p = \|\mathbf{Px}\|.$$

试证明 $\|\mathbf{x}\|_p$ 是 \mathbf{R}^n 上向量的一种范数.

14. 设 $\mathbf{A} \in \mathbf{R}^{n \times n}$ 为对称正定阵, 定义

$$\|\mathbf{x}\|_{\mathbf{A}} = (\mathbf{Ax}, \mathbf{x})^{\frac{1}{2}},$$

试证明 $\|\mathbf{x}\|_{\mathbf{A}}$ 为 \mathbf{R}^n 上向量的一种范数.

15. 设 $\|\mathbf{A}\|_s$, $\|\mathbf{A}\|_t$ 为 $\mathbf{R}^{n \times n}$ 上任意两种矩阵算子范数, 证明存在常数 $c_1, c_2 > 0$, 使对一切 $\mathbf{A} \in \mathbf{R}^{n \times n}$ 满足

$$c_1 \|\mathbf{A}\|_s \leq \|\mathbf{A}\|_t \leq c_2 \|\mathbf{A}\|_s.$$

16. 设 \mathbf{A} 为非奇异矩阵, 求证

$$\frac{1}{\mathbf{A}^{-1}} = \min_{\mathbf{y} \neq \mathbf{0}} \frac{\|\mathbf{Ay}\|_2}{\|\mathbf{y}\|_2}.$$

17. 矩阵第一行乘以一数, 成为

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

证明当 $= \pm \frac{2}{3}$ 时, $\text{cond}(\mathbf{A})$ 有最小值.

18. 设

$$\mathbf{A} = \begin{pmatrix} 100 & 99 \\ 99 & 98 \end{pmatrix},$$

计算 \mathbf{A} 的条件数 $\text{cond}(\mathbf{A})_v$ ($v=2, \dots$).

19. 证明: 如果 \mathbf{A} 是正交阵, 则 $\text{cond}(\mathbf{A})_2 = 1$.

20. 设 $\mathbf{A}, \mathbf{B} \in \mathbf{R}^{n \times n}$, 且 \cdot 为 $\mathbf{R}^{n \times n}$ 上矩阵的算子范数, 证明:

$$\text{cond}(\mathbf{AB}) = \text{cond}(\mathbf{A})\text{cond}(\mathbf{B}).$$

21. 设 $\mathbf{Ax} = \mathbf{b}$, 其中 $\mathbf{A} \in \mathbf{R}^{n \times n}$ 为非奇异阵, 证明:

(a) $\mathbf{A}^T \mathbf{A}$ 为对称正定矩阵;

(b) $\text{cond}(\mathbf{A}^T \mathbf{A})_2 = [\text{cond}(\mathbf{A})_2]^2$.

第 6 章 解线性方程组的迭代法

6.1 引言

考虑线性方程组

$$\mathbf{A}\mathbf{x} = \mathbf{b}, \quad (1.1)$$

其中 \mathbf{A} 为非奇异矩阵, 当 \mathbf{A} 为低阶稠密矩阵时, 第 5 章所讨论的选主元消去法是解(1.1)的有效方法. 但是, 对于由工程技术中产生的大型稀疏矩阵方程组 (\mathbf{A} 的阶数 n 很大, 但零元素较多, 例如求某些偏微分方程数值解所产生的线性方程组, $n \approx 10^4$), 利用迭代法求解(1.1)是合适的. 在计算机内存和运算两方面, 迭代法通常都可利用 \mathbf{A} 中有大量零元素的特点.

本章将介绍迭代法的一些基本理论及雅可比迭代法, 高斯-塞德尔迭代法, 超松弛迭代法, 而超松弛迭代法应用很广泛.

下面举简例, 以便了解迭代法的思想.

例 1 求解方程组

$$\begin{aligned} 8x_1 - 3x_2 + 2x_3 &= 20, \\ 4x_1 + 11x_2 - x_3 &= 33, \\ 6x_1 + 3x_2 + 12x_3 &= 36. \end{aligned} \quad (1.2)$$

记为 $\mathbf{Ax} = \mathbf{b}$, 其中

$$\mathbf{A} = \begin{matrix} 8 & -3 & 2 \\ 4 & 11 & -1 \\ 6 & 3 & 12 \end{matrix}, \quad \mathbf{x} = \begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix}, \quad \mathbf{b} = \begin{matrix} 20 \\ 33 \\ 36 \end{matrix}.$$

方程组的精确解是 $\mathbf{x}^* = (3, 2, 1)^T$. 现将(1.2)改写为

$$\begin{aligned}x_1 &= \frac{1}{8}(3x_2 - 2x_3 + 20), \\x_2 &= \frac{1}{11}(-4x_1 + x_3 + 33), \\x_3 &= \frac{1}{12}(-6x_1 - 3x_2 + 36).\end{aligned}\quad (1.3)$$

或写为 $\mathbf{x} = \mathbf{B} \mathbf{x} + \mathbf{f}$, 其中

$$\mathbf{B} = \begin{pmatrix} 0 & \frac{3}{8} & -\frac{2}{8} & \frac{20}{8} \\ -\frac{4}{11} & 0 & \frac{1}{11} & \frac{33}{11} \\ -\frac{6}{12} & -\frac{3}{12} & 0 & \frac{36}{12} \end{pmatrix}, \quad \mathbf{f} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}.$$

我们任取初始值, 例如取 $\mathbf{x}^{(0)} = (0, 0, 0)^T$. 将这些值代入 (1.3) 式右边 (若 (1.3) 式为等式即求得方程组的解, 但一般不满足), 得到新的值 $\mathbf{x}^{(1)} = (x_1^{(1)}, x_2^{(1)}, x_3^{(1)})^T = (2.5, 3, 3)^T$, 再将 $\mathbf{x}^{(1)}$ 分量代入 (1.3) 式右边得到 $\mathbf{x}^{(2)}$, 反复利用这个计算程序, 得到一向量序列和一般的计算公式 (迭代公式)

$$\begin{aligned}\mathbf{x}^{(0)} &= \begin{pmatrix} x_1^{(0)} \\ x_2^{(0)} \\ x_3^{(0)} \end{pmatrix}, \quad \mathbf{x}^{(1)} = \begin{pmatrix} x_1^{(1)} \\ x_2^{(1)} \\ x_3^{(1)} \end{pmatrix}, \quad \dots, \quad \mathbf{x}^{(k)} = \begin{pmatrix} x_1^{(k)} \\ x_2^{(k)} \\ x_3^{(k)} \end{pmatrix}, \quad \dots \\x_1^{(k+1)} &= (3x_2^{(k)} - 2x_3^{(k)} + 20)/8, \\x_2^{(k+1)} &= (-4x_1^{(k)} + x_3^{(k)} + 33)/11, \\x_3^{(k+1)} &= (-6x_1^{(k)} - 3x_2^{(k)} + 36)/12.\end{aligned}\quad (1.4)$$

简写为

$$\mathbf{x}^{(k+1)} = \mathbf{B} \mathbf{x}^{(k)} + \mathbf{f},$$

其中 k 表示迭代次数 ($k = 0, 1, 2, \dots$) .

迭代到第 10 次有

$$\begin{aligned}\mathbf{x}^{(10)} &= (3.000032, 1.999838, 0.9998813)^T; \\ \mathbf{x}^{(10)} &= 0.000187 (\mathbf{x}^{(10)} - \mathbf{x}^*) .\end{aligned}$$

从此例看出,由迭代法产生的向量序列 $\mathbf{x}^{(k)}$ 逐步逼近方程组的精确解 \mathbf{x}^* .

对于任何一个方程组 $\mathbf{x} = \mathbf{Bx} + \mathbf{f}$ (由 $\mathbf{Ax} = \mathbf{b}$ 变形得到的等价方程组),由迭代法产生的向量序列 $\mathbf{x}^{(k)}$ 是否一定逐步逼近方程组的解 \mathbf{x}^* 呢?回答是不一定.请读者考虑用迭代法解下述方程组

$$\begin{aligned}x_1 &= 2x_2 + 5, \\ x_2 &= 3x_1 + 5.\end{aligned}$$

对于给定方程组 $\mathbf{x} = \mathbf{Bx} + \mathbf{f}$,设有唯一解 \mathbf{x}^* ,则

$$\mathbf{x}^* = \mathbf{Bx}^* + \mathbf{f}. \quad (1.5)$$

又设 $\mathbf{x}^{(0)}$ 为任取的初始向量,按下述公式构造向量序列

$$\mathbf{x}^{(k+1)} = \mathbf{Bx}^{(k)} + \mathbf{f}, \quad k = 0, 1, 2, \dots, \quad (1.6)$$

其中 k 表迭代次数.

定义 1 (1) 对于给定的方程组 $\mathbf{x} = \mathbf{Bx} + \mathbf{f}$,用公式(1.6)逐步代入求近似解的方法称为迭代法(或称为一阶定常迭代法,这里 \mathbf{B} 与 k 无关).

(2) 如果 $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)}$ 存在(记为 \mathbf{x}^*),称此迭代法收敛,显然 \mathbf{x}^* 就是方程组的解,否则称此迭代法发散.

由上述讨论,需要研究 $\{\mathbf{x}^{(k)}\}$ 的收敛性.引进误差向量

$$\mathbf{e}^{(k+1)} = \mathbf{x}^{(k+1)} - \mathbf{x}^*,$$

由(1.6)减去(1.5)式,得 $\mathbf{e}^{(k+1)} = \mathbf{B} \mathbf{e}^{(k)} (k = 0, 1, 2, \dots)$,递推得

$$\mathbf{e}^{(k)} = \mathbf{B}^{(k-1)} = \dots = \mathbf{B}^{(0)} \mathbf{e}^{(0)}.$$

要考察 $\{\mathbf{x}^{(k)}\}$ 的收敛性.就要研究 \mathbf{B} 在什么条件下有 $\lim_{k \rightarrow \infty} \mathbf{B}^{(k)} = \mathbf{0}$,亦即要研究 \mathbf{B} 满足什么条件时有 $\mathbf{B}^k = \mathbf{0}$ (零矩阵)($k = 0, 1, 2, \dots$).

6.2 基本迭代法

设有

$$\mathbf{A}\mathbf{x} = \mathbf{b}, \quad (2.1)$$

其中, $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{n \times n}$ 为非奇异矩阵. 下面研究如何建立解 $\mathbf{A}\mathbf{x} = \mathbf{b}$ 的各种迭代法.

将 \mathbf{A} 分裂为

$$\mathbf{A} = \mathbf{M} - \mathbf{N}, \quad (2.2)$$

其中, \mathbf{M} 为可选择的非奇异矩阵, 且使 $\mathbf{M}\mathbf{x} = \mathbf{d}$ 容易求解, 一般选择为 \mathbf{A} 的某种近似, 称 \mathbf{M} 为分裂矩阵.

于是, 求解 $\mathbf{A}\mathbf{x} = \mathbf{b}$ 转化为求解 $\mathbf{M}\mathbf{x} = \mathbf{N}\mathbf{x} + \mathbf{b}$, 即求解

$$\mathbf{A}\mathbf{x} = \mathbf{b} \quad \text{求解 } \mathbf{x} = \mathbf{M}^{-1} \mathbf{N}\mathbf{x} + \mathbf{M}^{-1} \mathbf{b}.$$

可构造一阶定常迭代法

$$\begin{aligned} & \mathbf{x}^{(0)} \text{ (初始向量),} \\ & \mathbf{x}^{(k+1)} = \mathbf{B}\mathbf{x}^{(k)} + \mathbf{f} \quad (k = 0, 1, \dots,), \end{aligned} \quad (2.3)$$

其中 $\mathbf{B} = \mathbf{M}^{-1} \mathbf{N} = \mathbf{M}^{-1} (\mathbf{M} - \mathbf{A}) = \mathbf{I} - \mathbf{M}^{-1} \mathbf{A}$, $\mathbf{f} = \mathbf{M}^{-1} \mathbf{b}$. 称 $\mathbf{B} = \mathbf{I} - \mathbf{M}^{-1} \mathbf{A}$ 为迭代法的迭代矩阵, 选取 \mathbf{M} 阵, 就得到解 $\mathbf{A}\mathbf{x} = \mathbf{b}$ 的各种迭代法.

设 $a_{ii} \neq 0$ ($i = 1, 2, \dots, n$), 并将 \mathbf{A} 写为三部分

$$\mathbf{A} = \begin{matrix} & & 0 & & & & \\ a_{11} & & -a_{11} & & & & 0 \\ & a_{22} & & \cdots & & \cdots & & w \\ & & - & \cdots & & \cdots & & \\ w & & & -a_{n-1,1} & -a_{n-1,2} & \cdots & & 0 \\ a_{nn} & & & -a_{n1} & -a_{n2} & \cdots & -a_{n,n-1} & 0 \end{matrix}$$

$$\begin{array}{ccccccccc}
 & 0 & - & a_{12} & \cdots & - & a_{1,n-1} & - & a_{1n} \\
 & 0 & \cdots & - & a_{2,n-1} & - & a_{2n} \\
 & \vdots & & \vdots & \cdots & & \vdots & & \vdots \\
 & & & w & \cdots & & \cdots & & \\
 & & & 0 & & - & a_{n-1,n} \\
 & & & & & & 0
 \end{array}
 \quad \mathbf{D} = \mathbf{L} + \mathbf{U}. \quad (2.4)$$

6.2.1 雅可比迭代法

由 $a_{ii} \neq 0$ ($i = 1, \dots, n$), 选取 \mathbf{M} 为 \mathbf{A} 的对角元素部分, 即选取 $\mathbf{M} = \mathbf{D}$ (对角阵), $\mathbf{A} = \mathbf{D} - \mathbf{N}$, 由 (2.3) 式得到解 $\mathbf{Ax} = \mathbf{b}$ 的雅可比 (Jacobi) 迭代法

$$\begin{aligned}
 \mathbf{x}^{(0)} & \quad (\text{初始向量}) \\
 \mathbf{x}^{(k+1)} & = \mathbf{Bx}^{(k)} + \mathbf{f} \quad (k = 0, 1, \dots),
 \end{aligned} \quad (2.5)$$

其中 $\mathbf{B} = \mathbf{I} - \mathbf{D}^{-1} \mathbf{A} = \mathbf{D}^{-1} (\mathbf{L} + \mathbf{U}) = \mathbf{J}$, $\mathbf{f} = \mathbf{D}^{-1} \mathbf{b}$. 称 \mathbf{J} 为解 $\mathbf{Ax} = \mathbf{b}$ 的雅可比迭代法的迭代阵.

下面给出雅可比迭代法 (2.5) 的分量计算公式, 记

$$\mathbf{x}^{(k)} = (x_1^{(k)}, \dots, x_i^{(k)}, \dots, x_n^{(k)})^T,$$

由雅可比迭代公式 (2.5) 有

$$\mathbf{Dx}^{(k+1)} = (\mathbf{L} + \mathbf{U}) \mathbf{x}^{(k)} + \mathbf{b},$$

$$\text{或 } a_{ii} x_i^{(k+1)} = - \sum_{j=1}^{i-1} a_{ij} x_j^{(k)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} + b \quad (i = 1, 2, \dots, n).$$

于是, 解 $\mathbf{Ax} = \mathbf{b}$ 的雅可比迭代法的计算公式为

$$\begin{aligned}
 \mathbf{x}^{(0)} & = (x_1^{(0)}, \dots, x_n^{(0)})^T, \\
 x_i^{(k+1)} & = (b - \sum_{j=1}^{i-1} a_{ij} x_j^{(k)}) / a_{ii} \\
 (i & = 1, 2, \dots, n) \quad (k = 0, 1, \dots \text{ 表示迭代次数}).
 \end{aligned} \quad (2.6)$$

由(2.6)式可知, 雅可比迭代法计算公式简单, 每迭代一次只需计算一次矩阵和向量的乘法且计算过程中原始矩阵 \mathbf{A} 始终不变.

6.2.2 高斯-塞德尔迭代法

选取分裂矩阵 \mathbf{M} 为 \mathbf{A} 的下三角部分, 即选取 $\mathbf{M} = \mathbf{D} - \mathbf{L}$ (下三角阵), $\mathbf{A} = \mathbf{M} + \mathbf{N}$, 于是由(2.3)式得到解 $\mathbf{Ax} = \mathbf{b}$ 的高斯-塞德尔(Gauss-Seidel)迭代法

$$\begin{aligned} \mathbf{x}^{(0)} & \quad (\text{初始向量}) \\ \mathbf{x}^{(k+1)} &= \mathbf{Bx}^{(k)} + \mathbf{f} \quad (k = 0, 1, \dots), \end{aligned} \quad (2.7)$$

其中 $\mathbf{B} = \mathbf{I} - (\mathbf{D} - \mathbf{L})^{-1} \mathbf{A} = (\mathbf{D} - \mathbf{L})^{-1} \mathbf{U}$, \mathbf{G} , $\mathbf{f} = (\mathbf{D} - \mathbf{L})^{-1} \mathbf{b}$. 称 $\mathbf{G} = (\mathbf{D} - \mathbf{L})^{-1} \mathbf{U}$ 为解 $\mathbf{Ax} = \mathbf{b}$ 的高斯-塞德尔迭代法的迭代阵.

下面给出高斯-塞德尔迭代法的分量计算公式. 记

$$\mathbf{x}^{(k)} = (x_1^{(k)}, \dots, x_i^{(k)}, \dots, x_n^{(k)})^T.$$

由(2.7)式有

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

或

$$\mathbf{Dx}^{(k+1)} = \mathbf{Lx}^{(k+1)} + \mathbf{Ux}^{(k)} + \mathbf{b},$$

$$\text{即 } a_{ii} x_i^{(k+1)} = b - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} \quad (i = 1, 2, \dots, n).$$

于是解 $\mathbf{Ax} = \mathbf{b}$ 的高斯-塞德尔迭代法计算公式为

$$\begin{aligned} \mathbf{x}^{(0)} &= (x_1^{(0)}, \dots, x_n^{(0)})^T \quad (\text{初始向量}), \\ x_i^{(k+1)} &= b - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} \quad / a_{ii} \quad (2.8) \\ (i &= 1, \dots, n; k = 0, 1, \dots). \end{aligned}$$

或

$$\begin{aligned}
 \mathbf{x}^{(0)} &= (x_1^{(0)}, \dots, x_n^{(0)})^T, \\
 x_i^{(k+1)} &= x_i^{(k)} + x_i \\
 x_i &= b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} \quad | / a_{ii} \quad (2.9) \\
 (i &= 1, \dots, n; k = 0, 1, \dots).
 \end{aligned}$$

雅可比迭代法不使用变量的最新信息计算 $x_i^{(k+1)}$, 而由高斯-塞德尔迭代公式(2.8)可知, 计算 $\mathbf{x}^{(k+1)}$ 的第 i 个分量 $x_i^{(k+1)}$ 时, 利用了已经计算出的最新分量 $x_j^{(k+1)}$ ($j = 1, 2, \dots, i-1$)。高斯-塞德尔迭代法可看作雅可比迭代法的一种改进。由(2.8)可知, 高斯-塞德尔迭代法每迭代一次只需计算一次矩阵与向量的乘法。

算法 1(高斯-塞德尔迭代法) 设 $\mathbf{Ax} = \mathbf{b}$, 其中 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 为非奇异矩阵且 $a_{ii} \neq 0$ ($i = 1, \dots, n$), 本算法用高斯-塞德尔迭代法解 $\mathbf{Ax} = \mathbf{b}$, 数组 $x(n)$ 开始存放 $\mathbf{x}^{(0)}$, 后存放 $\mathbf{x}^{(k)}$, N_0 为最大迭代次数。

$$1. \quad x_i = 0.0 \quad (i = 1, \dots, n)$$

2. 对于 $k = 1, 2, \dots, N_0$

对于 $i = 1, 2, \dots, n$

$$x_i = b_i - \sum_{j=1}^{i-1} a_{ij} x_j - \sum_{j=i+1}^n a_{ij} x_j \quad | / a_{ii}$$

迭代一次, 这个算法需要的运算次数至多与矩阵 \mathbf{A} 的非零元素的个数一样多。

例 2 用高斯-塞德尔迭代法解线性方程组(1.2)。

高斯-塞德尔迭代公式: 取 $\mathbf{x}^{(0)} = (0, 0, 0)^T$ 。

$$\begin{aligned}
 x_1^{(k+1)} &= (20 + x_2^{(k)} - 2x_3^{(k)}) / 8, \\
 x_2^{(k+1)} &= (33 - 4x_1^{(k+1)} + x_3^{(k)}) / 11, \\
 x_3^{(k+1)} &= (36 - 6x_1^{(k+1)} - 3x_2^{(k+1)}) / 12.
 \end{aligned}$$

$$(k = 0, 1, \dots,)$$

计算 $\mathbf{x}^{(7)} = (3.000002, 1.9999987, 0.9999932)^T$, 且

$$\|\mathbf{x}^* - \mathbf{x}^{(7)}\| < 2.02 \times 10^{-6}.$$

由此例可知, 用高斯-塞德尔迭代法, 雅可比迭代法解线性方程组(1.2)(且取 $\mathbf{x}^{(0)} = \mathbf{0}$)均收敛, 而高斯-塞德尔迭代法比雅可比迭代法收敛较快(即取 $\mathbf{x}^{(0)}$ 相同, 达到同样精度所需迭代次数较少), 但这结论只当 \mathbf{A} 满足一定条件时才是对的.

6.2.3 解大型稀疏线性方程组的逐次超松弛迭代法

选取分裂矩阵 \mathbf{M} 为带参数的下三角阵

$$\mathbf{M} = \frac{1}{\omega} (\mathbf{D} - \mathbf{L}),$$

其中 $\omega > 0$ 为可选择的松弛因子.

于是, 由(2.3)可构造一个迭代法, 其迭代矩阵为

$$\begin{aligned} \mathbf{L} &= \mathbf{I} - (\mathbf{D} - \mathbf{L})^{-1} \mathbf{A} \\ &= (\mathbf{D} - \mathbf{L})^{-1} ((1 - \omega) \mathbf{D} + \omega \mathbf{U}). \end{aligned}$$

从而得到解 $\mathbf{Ax} = \mathbf{b}$ 的逐次超松弛迭代法 (Successive Over Relaxation Method, 简称 SOR 方法).

解 $\mathbf{Ax} = \mathbf{b}$ 的 SOR 方法为

$$\begin{aligned} \mathbf{x}^{(0)} &\quad (\text{初始向量}) \\ \mathbf{x}^{(k+1)} &= \mathbf{L} \mathbf{x}^{(k)} + \mathbf{f} \quad (k = 0, 1, \dots), \end{aligned} \tag{2.10}$$

其中 $\mathbf{L} = (\mathbf{D} - \mathbf{L})^{-1} ((1 - \omega) \mathbf{D} + \omega \mathbf{U})$, $\mathbf{f} = (\mathbf{D} - \mathbf{L})^{-1} \mathbf{b}$.

下面给出解 $\mathbf{Ax} = \mathbf{b}$ 的 SOR 迭代法的分量计算公式. 记

$$\mathbf{x}^{(k)} = (x_1^{(k)}, \dots, x_i^{(k)}, \dots, x_n^{(k)})^T,$$

由(2.10)式可得

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

或

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

由此, 得到解 $\mathbf{Ax} = \mathbf{b}$ 的 SOR 方法的计算公式

$$\begin{aligned} \mathbf{x}^{(0)} &= (x_1^{(0)}, \dots, x_n^{(0)})^T, \\ x_i^{(k+1)} &= x_i^{(k)} + b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i}^n a_{ij} x_j^{(k)} / a_{ii} \\ (i &= 1, 2, \dots, n; k = 0, 1), \\ \text{为松弛因子.} \end{aligned} \tag{2.11}$$

或

$$\begin{aligned} \mathbf{x}^{(0)} &= (x_1^{(0)}, \dots, x_n^{(0)})^T, \\ x_i^{(k+1)} &= x_i^{(k)} + x_i, \\ x_i &= b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i}^n a_{ij} x_j^{(k)} / a_{ii} \\ (i &= 1, 2, \dots, n; k = 0, 1, \dots), \\ \text{为松弛因子.} \end{aligned} \tag{2.12}$$

(1) 显然, 当 $\omega = 1$ 时, SOR 方法即为高斯-塞德尔迭代法.

(2) SOR 方法每迭代一次主要运算量是计算一次矩阵与向量的乘法.

(3) 当 $\omega > 1$ 时, 称为超松弛法; 当 $\omega < 1$ 时, 称为低松弛法.

(4) 在计算机实现时可用

$$\max_{1 \leq i \leq n} |x_i| = \max_{1 \leq i \leq n} |x_i^{(k+1)} - x_i^{(k)}| <$$

控制迭代终止, 或用 $\mathbf{r}^{(k)} = \mathbf{b} - \mathbf{Ax}^{(k)} <$ 控制迭代终止.

SOR 迭代法是高斯-塞德尔迭代法的一种修正, 可由下述思想得到.

设已知 $\mathbf{x}^{(k)}$ 及已计算 $\mathbf{x}^{(k+1)}$ 的分量 $x_j^{(k+1)} (j = 1, 2, \dots, i-1)$.

(1) 首先用高斯-塞德尔迭代法定义辅助量 $\bar{x}_i^{(k+1)}$,

$$\bar{x}_i^{(k+1)} = b - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} / a_{ii}. \quad (2.13)$$

(2) 再由 $x_i^{(k)}$ 与 $\bar{x}_i^{(k+1)}$ 加权平均定义 $x_i^{(k+1)}$, 即

$$x_i^{(k+1)} = (1 - \omega) x_i^{(k)} + \bar{x}_i^{(k+1)} = x_i^{(k)} + (\bar{x}_i^{(k+1)} - x_i^{(k)}) . \quad (2.14)$$

将(2.13)代入(2.14)得到解 $\mathbf{Ax} = \mathbf{b}$ 的 SOR 迭代(2.11)式.

例 3 用 SOR 方法解方程组

$$\begin{array}{cccccc} -4 & 1 & 1 & 1 & x_1 & 1 \\ 1 & -4 & 1 & 1 & x_2 & 1 \\ 1 & 1 & -4 & 1 & x_3 & 1 \\ 1 & 1 & 1 & -4 & x_4 & 1 \end{array} ,$$

它的精确解为 $\mathbf{x}^* = (-1, -1, -1, -1)^T$.

解 取 $\mathbf{x}^{(0)} = \mathbf{0}$, 迭代公式为

$$\begin{aligned} x_1^{(k+1)} &= x_1^{(k)} - (1 + 4 x_1^{(k)} - x_2^{(k)} - x_3^{(k)} - x_4^{(k)}) / 4; \\ x_2^{(k+1)} &= x_2^{(k)} - (1 - x_1^{(k+1)} + 4 x_2^{(k)} - x_3^{(k)} - x_4^{(k)}) / 4; \\ x_3^{(k+1)} &= x_3^{(k)} - (1 - x_1^{(k+1)} - x_2^{(k+1)} + 4 x_3^{(k)} - x_4^{(k)}) / 4; \\ x_4^{(k+1)} &= x_4^{(k)} - (1 - x_1^{(k+1)} - x_2^{(k+1)} - x_3^{(k+1)} + 4 x_4^{(k)}) / 4 . \end{aligned}$$

取 $\omega = 1.3$, 第 11 次迭代结果为

$$\begin{aligned} \mathbf{x}^{(11)} &= (-0.99999646, -1.00000310, -0.99999953, \\ &\quad -0.99999912)^T, \\ &\quad (11) \quad 2 \quad 0.46 \times 10^{-5} . \end{aligned}$$

对 ω 取其他值, 迭代次数如下表. 从此例看到, 松弛因子选择得好, 会使 SOR 迭代法的收敛大大加速. 本例中 $\omega = 1.3$ 是最佳松弛因子.

表 6-1

松弛因子	满足误差 $\ \mathbf{x}^{(k)} - \mathbf{x}^*\ _2 < 10^{-5}$ 的迭代次数	松弛因子	满足误差 $\ \mathbf{x}^{(k)} - \mathbf{x}^*\ _2 < 10^{-5}$ 的迭代次数
1.0	22	1.5	17
1.1	17	1.6	23
1.2	12	1.7	33
1.3	11 (最少迭代次数)	1.8	53
1.4	14	1.9	109

6.3 迭代法的收敛性

6.3.1 一阶定常迭代法的基本定理

设

$$\mathbf{Ax} = \mathbf{b}, \quad (3.1)$$

其中 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 为非奇异矩阵, 记 \mathbf{x}^* 为(3.1)精确解, 且设有等价的方程组

$$\mathbf{Ax} = \mathbf{b} \quad \mathbf{x} = \mathbf{Bx} + \mathbf{f}.$$

于是

$$\mathbf{x}^* = \mathbf{Bx}^* + \mathbf{f}. \quad (3.2)$$

设有解 $\mathbf{Ax} = \mathbf{b}$ 的一阶定常迭代法

$$\mathbf{x}^{(k+1)} = \mathbf{Bx}^{(k)} + \mathbf{f}. \quad (3.3)$$

有意义的问题是: 迭代矩阵 \mathbf{B} 满足什么条件时, 由迭代法产生的向量序列 $\{\mathbf{x}^{(k)}\}$ 收敛到 \mathbf{x}^* .

引进误差向量

$$\mathbf{e}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}^* \quad (k = 0, 1, 2, \dots).$$

由(3.3)式减(3.2)式得到误差向量的递推公式

$$\begin{aligned} & {}^{(k+1)} = \mathbf{B}^{(k)}, \\ & {}^{(k)} = \mathbf{B}^{(0)} \quad (k = 0, 1, \dots). \end{aligned}$$

由 6.1 节可知, 研究迭代法(3.3)收敛性问题就是要研究迭代矩阵 \mathbf{B} 满足什么条件时, 有 $\mathbf{B}^k \rightarrow \mathbf{0}$ (零矩阵) ($k \rightarrow \infty$) .

定义 2 设有矩阵序列 $\mathbf{A}_k = (a_{ij}^{(k)}) \in \mathbb{R}^{n \times n}$ 及 $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{n \times n}$,

如果 n^2 个数列极限存在且有

$$\lim_k a_{ij}^{(k)} = a_{ij} \quad (i, j = 1, 2, \dots, n),$$

则称 $\{\mathbf{A}_k\}$ 收敛于 \mathbf{A} , 记为 $\lim_k \mathbf{A}_k = \mathbf{A}$.

例 4 设有矩阵序列

$$\mathbf{A} = \begin{pmatrix} 1 & & & \\ 0 & 2 & & \\ & & 2 & \\ & & & \ddots \end{pmatrix}, \quad \mathbf{A}^2 = \begin{pmatrix} 2 & & & \\ 0 & 2 & & \\ & & 2 & \\ & & & \ddots \end{pmatrix}, \quad \dots, \quad \mathbf{A}^k = \begin{pmatrix} k & & & \\ 0 & k & & \\ & & k & \\ & & & \ddots \end{pmatrix}, \quad \dots$$

且设 $|k| < 1$, 考查其极限.

解 显然, 当 $|k| < 1$ 时, 则有 $\lim_k \mathbf{A}_k = \lim_k \mathbf{A}^k = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$.

矩阵序列极限概念可以用矩阵算子范数来描述.

定理 1 $\lim_k \mathbf{A}_k = \mathbf{A} \quad \lim_k \|\mathbf{A}_k - \mathbf{A}\| = 0$, 其中 $\|\cdot\|$ 为矩阵的任意一种算子范数.

证明 显然有

$$\lim_k \mathbf{A}_k = \mathbf{A} \quad \lim_k \|\mathbf{A}_k - \mathbf{A}\| = 0.$$

再利用矩阵范数的等价性, 可证定理对其他算子范数亦对.

定理 2 $\lim_k \mathbf{A}_k = \mathbf{A}$ 是对任何向量 $\mathbf{x} \in \mathbb{R}^n$ 都有 $\lim_k \mathbf{A}_k \mathbf{x} = \mathbf{A} \mathbf{x}$.

证明作为练习.

定理 3 设 $\mathbf{B} = (b_{ij}) \in \mathbb{R}^{n \times n}$, 则 $\lim_k \mathbf{B}^k = \mathbf{0}$ (零矩阵) 的充分必要条件是矩阵 \mathbf{B} 的谱半径 $(\mathbf{B}) < 1$.

证明 由矩阵 \mathbf{B} 的若当标准型, 存在非奇异矩阵 \mathbf{P} 使

$$\mathbf{P}^{-1} \mathbf{B} \mathbf{P} = \begin{matrix} \mathbf{J} \\ \mathbf{W} \end{matrix} \quad \mathbf{J}_r$$

其中若当块

$$\mathbf{J}_i = \begin{matrix} i & 1 \\ \mathbf{W} & \\ \mathbf{W} & 1 \end{matrix},$$

$$i = n_i \times n_i$$

且 $\sum_{i=1}^r n_i = n$, 显然有

$$\mathbf{B} = \mathbf{P} \mathbf{J} \mathbf{P}^{-1},$$

$$\mathbf{B}^k = \mathbf{P} \mathbf{J}^k \mathbf{P}^{-1},$$

其中

$$\mathbf{J}^k = \begin{matrix} \mathbf{J} \\ \mathbf{W} \\ \mathbf{W} \\ \vdots \\ \mathbf{J}_r \end{matrix}.$$

于是 $\lim_k \mathbf{B}^k = \mathbf{0}$ $\lim_k \mathbf{J}^k = \mathbf{0}$ $\lim_k \mathbf{J}_i^k = \mathbf{0}$ ($i = 1, 2, \dots, r$) .

下面考查 \mathbf{J}_i^k 的情况. 引进记号

$$\mathbf{E}_{i,k} = \begin{matrix} \mathbf{0} & \mathbf{I} & t - k \\ \mathbf{0} & \mathbf{0} \end{matrix} \quad \mathbf{R}^{t \times t} \quad (t = n_i).$$

显然有, $\mathbf{E}_{i,0} = \mathbf{I}$, $\mathbf{E}_{i,k} = \mathbf{0}$ (当 $k > t$), $(\mathbf{E}_{i,1})^k = \mathbf{E}_{i,k}$. 由于 $\mathbf{J}_i =$
 $i \mathbf{I} + \mathbf{E}_{i,1}$, 因此

$$\mathbf{J}_i^k = (i \mathbf{I} + \mathbf{E}_{i,1})^k = \sum_{j=0}^k \mathbf{C}_k^j \begin{matrix} k-j \\ i \end{matrix} (\mathbf{E}_{i,1})^j$$

$$\begin{aligned}
 &= \sum_{j=0}^{t-1} C_k^{j-k-j} E_{t,j} \\
 &\quad C_k^1 \quad C_k^2 \quad \dots \quad C_k^{t-1} \\
 &\quad \vdots \quad \vdots \quad \dots \quad \vdots \\
 &= \begin{matrix} W & W & \dots \\ W & C_k^1 \quad \vdots \quad \vdots \quad \vdots \\ & \vdots & \ddots & \vdots \\ & & & t \times t \end{matrix} \\
 &\quad (i = 1, 2, \dots, r),
 \end{aligned}$$

其中

$$C_k^j = \frac{k!}{j!(k-j)!} = \frac{k(k-1)\dots(k-j+1)}{j!}.$$

利用极限 $\lim_k k^r c^k = 0 (0 < c < 1, r > 0)$, 得到

$$\lim_k C_k^{j-k-j} = 0 \quad / \quad / < 1.$$

所以 $\lim_k B_i = \mathbf{0}$ 的充要条件是 $|i| < 1 (i = 1, 2, \dots, r)$, 即 $(B) < 1$.

定理 4(迭代法基本定理) 设有方程组

$$\mathbf{x} = B\mathbf{x} + \mathbf{f}, \quad (3.4)$$

及一阶定常迭代法

$$\mathbf{x}^{(k+1)} = B\mathbf{x}^{(k)} + \mathbf{f}. \quad (3.5)$$

对任意选取初始向量 $\mathbf{x}^{(0)}$, 迭代法(3.5)收敛的充要条件是矩阵 B 的谱半径 $(B) < 1$.

证明 充分性. 设 $(B) < 1$, 易知 $A\mathbf{x} = \mathbf{f}$ (其中 $A = I - B$) 有唯一解, 记为 $\hat{\mathbf{x}}$, 则

$$\hat{\mathbf{x}} = B\hat{\mathbf{x}} + \mathbf{f},$$

误差向量

$$\mathbf{x}^{(k)} = \mathbf{x}^{(k)} - \hat{\mathbf{x}} = B^k \mathbf{x}^{(0)}, \quad \mathbf{x}^{(0)} = \mathbf{x}^{(0)} - \hat{\mathbf{x}}.$$

由设 $(B) < 1$, 应用定理 3, 有 $\lim_k B^k = \mathbf{0}$. 于是对任意 $\mathbf{x}^{(0)}$ 有 $\lim_k \mathbf{x}^{(k)}$

$= \mathbf{0}$, 即 $\lim_k \mathbf{x}^{(k)} = \mathbf{x}^*$.

必要性. 设对任意 $\mathbf{x}^{(0)}$ 有

$$\lim_k \mathbf{x}^{(k)} = \mathbf{x}^*,$$

其中 $\mathbf{x}^{(k+1)} = \mathbf{B}\mathbf{x}^{(k)} + \mathbf{f}$. 显然, 极限 \mathbf{x}^* 是方程组(3.4)的解, 且对任意 $\mathbf{x}^{(0)}$ 有

$$^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}^* = \mathbf{B}^{(0)} - \mathbf{0} \quad (k).$$

由定理2知

$$\lim_k \mathbf{B}^k = \mathbf{0},$$

再由定理3, 即得 $(\mathbf{B}) < 1$.

定理4是一阶定常迭代法的基本理论.

推论 设 $\mathbf{Ax} = \mathbf{b}$, 其中 $\mathbf{A} = \mathbf{D} - \mathbf{L} - \mathbf{U}$ 为非奇异矩阵且 \mathbf{D} 非奇异, 则

(1) 解方程组的雅可比迭代法收敛的充要条件是 $(\mathbf{J}) < 1$, 其中 $\mathbf{J} = \mathbf{D}^{-1}(\mathbf{L} + \mathbf{U})$.

(2) 解方程组的高斯-塞德尔迭代法收敛的充要条件是 $(\mathbf{G}) < 1$, 其中 $\mathbf{G} = (\mathbf{D} - \mathbf{L})^{-1}\mathbf{U}$.

(3) 解方程组的SOR方法收敛的充要条件是 $(\mathbf{L}) < 1$, 其中 $\mathbf{L} = (\mathbf{D} - \mathbf{L})^{-1}((1 - \gamma)\mathbf{D} + \mathbf{U})$.

例5 考察用雅可比方法解方程组(1.2)的收敛性. 迭代矩阵 \mathbf{J} 的特征方程为

$$\det(\mathbf{I} - \mathbf{J}) = -3 + 0.034090909 + 0.039772727 = 0,$$

解得 $\lambda_1 = -0.3082$,

$$\lambda_2 = 0.1541 + i0.3245,$$

$$\lambda_3 = 0.1541 - i0.3245,$$

$$|\lambda_2| = |\lambda_3| = 0.3592 < 1, |\lambda_1| < 1,$$

即 $(\mathbf{J}) < 1$. 所以用雅可比迭代法解方程组(1.2)是收敛的.

例 6 考察用迭代法解方程组

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

的收敛性, 其中 $\mathbf{B} = \begin{pmatrix} 0 & 2 \\ 3 & 0 \end{pmatrix}$, $\mathbf{f} = \begin{pmatrix} 5 \\ 5 \end{pmatrix}$.

解 特征方程为 $\det(\mathbf{I} - \mathbf{B}) = -6 = 0$, 特征根 $\lambda_{1,2} = \pm 6$, 即 $(\mathbf{B}) > 1$. 这说明用迭代法解此方程组不收敛.

迭代法的基本定理在理论上是重要的, 由于 $(\mathbf{B}) > 1$, 下面利用矩阵 \mathbf{B} 的范数建立判别迭代法收敛的充分条件.

定理 5(迭代法收敛的充分条件) 设有方程组

$$\mathbf{x} = \mathbf{Bx} + \mathbf{f}, \quad \mathbf{B} \in \mathbb{R}^{n \times n},$$

及一阶定常迭代法

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

如果有 \mathbf{B} 的某种算子范数 $\|\mathbf{B}\| = q < 1$, 则

(1) 迭代法收敛, 即对任取 $\mathbf{x}^{(0)}$ 有

$$\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^* \quad \text{且} \quad \mathbf{x}^* = \mathbf{Bx}^* + \mathbf{f}.$$

$$(2) \quad \mathbf{x}^* - \mathbf{x}^{(k)} = q^k (\mathbf{x}^* - \mathbf{x}^{(0)}).$$

$$(3) \quad \mathbf{x}^* - \mathbf{x}^{(k)} = \frac{q}{1-q} (\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}).$$

$$(4) \quad \mathbf{x}^* - \mathbf{x}^{(k)} = \frac{q^k}{1-q} (\mathbf{x}^{(1)} - \mathbf{x}^{(0)}).$$

证明 (1) 由基本定理 4 结论(1)是显然的.

(2) 显然有关系式 $\mathbf{x}^* - \mathbf{x}^{(k+1)} = \mathbf{B}(\mathbf{x}^* - \mathbf{x}^{(k)})$ 及

$$\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} = \mathbf{B}(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}).$$

于是有 (a) $\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} = q(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})$;

$$(b) \quad \mathbf{x}^* - \mathbf{x}^{(k+1)} = q(\mathbf{x}^* - \mathbf{x}^{(k)}).$$

反复利用(b)即得(2).

(3) 考查

$$\begin{aligned} \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} &= \mathbf{x}^* - \mathbf{x}^{(k)} - (\mathbf{x}^* - \mathbf{x}^{(k+1)}) \\ &\quad \mathbf{x}^* - \mathbf{x}^{(k)} - \mathbf{x}^* - \mathbf{x}^{(k+1)} \\ &\quad (1 - q) \mathbf{x}^* - \mathbf{x}^{(k)}, \end{aligned}$$

即 $\mathbf{x}^* - \mathbf{x}^{(k)} = \frac{1}{1 - q} \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}$

$$\frac{q}{1 - q} \mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}.$$

(4) 反复利用(a), 则得到(4).

6.3.2 关于解某些特殊方程组迭代法的收敛性

在科学及工程计算中, 要求解方程组 $\mathbf{Ax} = \mathbf{b}$, 其矩阵 \mathbf{A} 常常具有某些特性. 例如, \mathbf{A} 具有对角占优性质或 \mathbf{A} 为不可约阵, 或 \mathbf{A} 是对称正定阵等, 下面讨论用基本迭代法解这些方程组的收敛性.

定义 3(对角占优阵) 设 $\mathbf{A} = (a_{ij})_{n \times n}$.

(1) 如果 \mathbf{A} 的元素满足

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

称 \mathbf{A} 为严格对角占优阵.

(2) 如果 \mathbf{A} 元素满足

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

且上式至少有一个不等式严格成立, 称 \mathbf{A} 为弱对角占优阵.

定义 4(可约与不可约矩阵) 设 $\mathbf{A} = (a_{ij})_{n \times n}$ ($n \geq 2$), 如果存在置换阵 \mathbf{P} 使

$$\mathbf{P}^T \mathbf{AP} = \begin{pmatrix} \mathbf{A}_1 & \mathbf{A}_2 \\ \mathbf{0} & \mathbf{A}_2 \end{pmatrix}, \quad (3.6)$$

其中 \mathbf{A}_1 为 r 阶方阵, \mathbf{A}_2 为 $n - r$ 阶方阵 ($1 \leq r < n$), 则称 \mathbf{A} 为可约

矩阵, 否则, 如果不存在这样置换阵 \mathbf{P} 使 (3.6) 式成立, 则称 \mathbf{A} 为不可约矩阵.

\mathbf{A} 为可约矩阵意即 \mathbf{A} 可经过若干行列重排化为 (3.6) 或 $\mathbf{Ax} = \mathbf{b}$ 可化为两个低阶方程组求解(如果 \mathbf{A} 经过两行交换的同时进行相应两列的交换, 称对 \mathbf{A} 进行一次行列重排).

事实上, 由 $\mathbf{Ax} = \mathbf{b}$ 可化为

$$\mathbf{P}^T \mathbf{AP}(\mathbf{P}^T \mathbf{x}) = \mathbf{P}^T \mathbf{b},$$

且记 $\mathbf{y} = \mathbf{P}^T \mathbf{x} = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}$, $\mathbf{P}^T \mathbf{b} = \begin{pmatrix} \mathbf{d} \\ \mathbf{d} \end{pmatrix}$, 其中 \mathbf{y}, \mathbf{d} 为 r 维向量. 于是, 求解 $\mathbf{Ax} = \mathbf{b}$ 化为求解

$$\mathbf{A}_{11} \mathbf{y}_1 + \mathbf{A}_{12} \mathbf{y}_2 = \mathbf{d},$$

$$\mathbf{A}_{21} \mathbf{y}_1 + \mathbf{A}_{22} \mathbf{y}_2 = \mathbf{d}.$$

由上式第 2 个方程组求出 \mathbf{y}_2 , 再代入第 1 个方程组求出 \mathbf{y}_1 .

显然, 如果 \mathbf{A} 所有元素都非零, 则 \mathbf{A} 为不可约阵.

例 7 设有矩阵

$$\mathbf{A} = \begin{matrix} b & c \\ a & b & c \\ & w & w & w \\ & a_{n-1} & b_{n-1} & c_{n-1} \\ & a_n & b_n \end{matrix}, \quad \mathbf{B} = \begin{matrix} 4 & -1 & -1 & 0 \\ -1 & 4 & 0 & -1 \\ -1 & 0 & 4 & -1 \\ 0 & -1 & -1 & 4 \end{matrix}.$$

(a_i, b_i, c_i 都不为零)

则 \mathbf{A}, \mathbf{B} 都是不可约矩阵.

定理 6(对角占优定理) 如果 $\mathbf{A} = (a_{ij})_{n \times n}$ 为严格对角占优矩阵或 \mathbf{A} 为不可约弱对角占优矩阵, 则 \mathbf{A} 为非奇异矩阵.

证明 只就 \mathbf{A} 为严格对角占优阵证明此定理. 采用反证法, 如果 $\det(\mathbf{A}) = 0$, 则 $\mathbf{Ax} = \mathbf{0}$ 有非零解, 记为 $\mathbf{x} = (x_1, \dots, x_n)^T$, 则 $|x_k| = \max_{\substack{i \\ 1 \leq i \leq n}} |x_i| \neq 0$.

由齐次方程组第 k 个方程

$$\sum_{j=1}^n a_{kj} x_j = 0,$$

则有

$$|a_{kk} x_k| = \left| \begin{array}{c|ccccc|c} & & & & & & \\ & a_{k1} & a_{k2} & \cdots & a_{kn} & & \\ \hline j=1 & & & & & & \\ j=k & & & & & & \\ \hline & a_{11} & a_{12} & \cdots & a_{1n} & & \\ & a_{21} & a_{22} & \cdots & a_{2n} & & \\ & \vdots & \vdots & \ddots & \vdots & & \\ & a_{n1} & a_{n2} & \cdots & a_{nn} & & \end{array} \right| x_k = |a_{kk}| |x_k|,$$

即

$$|a_{kk}| |x_k| = \left| \begin{array}{c|ccccc|c} & & & & & & \\ & a_{k1} & a_{k2} & \cdots & a_{kn} & & \\ \hline j=1 & & & & & & \\ j=k & & & & & & \\ \hline & a_{11} & a_{12} & \cdots & a_{1n} & & \\ & a_{21} & a_{22} & \cdots & a_{2n} & & \\ & \vdots & \vdots & \ddots & \vdots & & \\ & a_{n1} & a_{n2} & \cdots & a_{nn} & & \end{array} \right|,$$

与假设矛盾, 故 $\det(\mathbf{A}) = 0$.

定理 7 设 $\mathbf{Ax} = \mathbf{b}$, 如果:

(1) \mathbf{A} 为严格对角占优阵, 则解 $\mathbf{Ax} = \mathbf{b}$ 的雅可比迭代法, 高斯-塞德尔迭代法均收敛.

(2) \mathbf{A} 为弱对角占优阵, 且 \mathbf{A} 为不可约矩阵, 则解 $\mathbf{Ax} = \mathbf{b}$ 的雅可比迭代法, 高斯-塞德尔迭代法均收敛.

证明 只证(1)中高斯-塞德尔迭代法收敛, 其他同理可证.

由设可知, $a_{ii} \neq 0$ ($i = 1, \dots, n$), 解 $\mathbf{Ax} = \mathbf{b}$ 的高斯-塞德尔迭代法的迭代矩阵为 $\mathbf{G} = (\mathbf{D} - \mathbf{L})^{-1} \mathbf{U}$ ($\mathbf{A} = \mathbf{D} - \mathbf{L} - \mathbf{U}$). 下面考查 \mathbf{G} 的特征值情况.

$$\begin{aligned} \det(\mathbf{I} - \mathbf{G}) &= \det(\mathbf{I} - (\mathbf{D} - \mathbf{L})^{-1} \mathbf{U}) \\ &= \det((\mathbf{D} - \mathbf{L})^{-1}) \cdot \det((\mathbf{D} - \mathbf{L}) - \mathbf{U}). \end{aligned}$$

由于 $\det((\mathbf{D} - \mathbf{L})^{-1}) \neq 0$, 于是 \mathbf{G} 特征值即为 $\det((\mathbf{D} - \mathbf{L}) - \mathbf{U}) = 0$ 之根. 记

$$\mathbf{C} = (\mathbf{D} - \mathbf{L}) - \mathbf{U} = \begin{matrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \cdots & \cdots & \cdots & \cdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{matrix},$$

下面来证明, 当 $| \lambda | < 1$ 时, 则 $\det(\mathbf{C}) \neq 0$, 即 \mathbf{G} 的特征值均满足 $| \lambda | < 1$.

$\lambda < 1$, 由基本定理, 则有高斯-塞德尔迭代法收敛.

事实上, 当 $|\lambda| < 1$ 时, 由 \mathbf{A} 为严格对角占优阵, 则有

$$\begin{aligned} |c_{ii}| &= |a_{ii}| > \left| \sum_{j=1}^{i-1} |a_{ij}| + \sum_{j=i+1}^n |a_{ij}| \right| \\ &\quad \left| a_{ij} \right| + \left| a_{ij} \right| = \left| c_{ij} \right| \\ (i &= 1, 2, \dots, n). \end{aligned}$$

这说明, 当 $|\lambda| < 1$ 时, 矩阵 \mathbf{C} 为严格对角占优阵, 再由对角占优定理有 $\det(\mathbf{C}) \neq 0$.

下面研究对于解方程组 $\mathbf{Ax} = \mathbf{b}$ 的 SOR 方法中松弛因子 ω 在什么范围内取值, SOR 方法才可能收敛.

定理 8(SOR 方法收敛的必要条件) 设解方程组 $\mathbf{Ax} = \mathbf{b}$ 的 SOR 迭代法收敛, 则 $0 < \omega < 2$.

证明 由设 SOR 迭代法收敛, 则由定理 4 的推论中的(3)有 $(\mathbf{L}) < 1$, 设 \mathbf{L} 的特征值为 $\lambda_1, \lambda_2, \dots, \lambda_n$, 则

$$|\det(\mathbf{L})| = |\lambda_1 \lambda_2 \dots \lambda_n| = [(\mathbf{L})]^n,$$

或 $|\det(\mathbf{L})|^{1/n} = (\mathbf{L})^{1/n} < 1$.

另一方面

$$\begin{aligned} \det(\mathbf{L}) &= \det[(\mathbf{D} - \mathbf{L})^{-1}] \det((1 - \omega)\mathbf{D} + \omega\mathbf{U}) \\ &= (1 - \omega)^n, \end{aligned}$$

从而 $|\det(\mathbf{L})|^{1/n} = |1 - \omega|^{1/n} < 1$,

即 $0 < \omega < 2$.

定理 8 说明解 $\mathbf{Ax} = \mathbf{b}$ 的 SOR 迭代法, 只有在 $(0, 2)$ 范围内取松弛因子 ω , 才可能收敛.

定理 9 设 $\mathbf{Ax} = \mathbf{b}$, 如果:

(1) \mathbf{A} 为对称正定矩阵, $\mathbf{A} = \mathbf{D} - \mathbf{L} - \mathbf{U}$;

(2) $0 < \omega < 2$.

则解 $\mathbf{Ax} = \mathbf{b}$ 的 SOR 迭代法收敛.

证明 在上述假定下, 若能证明 $| \lambda | < 1$, 那么定理得证(其中为 \mathbf{L} 的任一特征值).

事实上, 设 \mathbf{y} 为对应的 \mathbf{L} 的特征向量, 即

$$\begin{aligned}\mathbf{L} \mathbf{y} &= \mathbf{y}, \quad \mathbf{y} = (y_1, y_2, \dots, y_n)^T - \mathbf{0}, \\ (\mathbf{D} - \mathbf{L})^{-1}((1 - \lambda) \mathbf{D} + \mathbf{U}) \mathbf{y} &= \mathbf{y},\end{aligned}$$

亦即

$$((1 - \lambda) \mathbf{D} + \mathbf{U}) \mathbf{y} = (\mathbf{D} - \mathbf{L}) \mathbf{y}.$$

为了找出 λ 的表达式, 考虑数量积

$$(((1 - \lambda) \mathbf{D} + \mathbf{U}) \mathbf{y}, \mathbf{y}) = ((\mathbf{D} - \mathbf{L}) \mathbf{y}, \mathbf{y}),$$

则

$$= \frac{(\mathbf{Dy}, \mathbf{y}) - (\mathbf{Dy}, \mathbf{y}) + (\mathbf{Uy}, \mathbf{y})}{(\mathbf{Dy}, \mathbf{y}) - (\mathbf{Ly}, \mathbf{y})},$$

显然

$$(\mathbf{Dy}, \mathbf{y}) = \sum_{i=1}^n a_{ii} / |y_i|^2 > 0, \quad (3.7)$$

记

$$- (\mathbf{Ly}, \mathbf{y}) = + i,$$

由于 $\mathbf{A} = \mathbf{A}^T$, 所以 $\mathbf{U} = \mathbf{L}^T$, 故

$$\begin{aligned}- (\mathbf{Uy}, \mathbf{y}) &= - (\mathbf{y}, \mathbf{Ly}) = - (\overline{\mathbf{Ly}}, \mathbf{y}) = - i, \\ 0 < (\mathbf{Ay}, \mathbf{y}) &= ((\mathbf{D} - \mathbf{L} - \mathbf{U}) \mathbf{y}, \mathbf{y}) = + 2,\end{aligned} \quad (3.8)$$

所以

$$= \frac{(- -) + i}{(+) + i},$$

从而

$$/ |^2 = \frac{(- -)^2 + ^2}{(+)^2 + ^2}.$$

当 $0 < < 2$ 时, 利用 (3.7), (3.8), 有

$$(-\lambda_1 - \lambda_2)^2 - (\lambda_1 + \lambda_2)^2 = (-2)(+2) < 0,$$

即 \mathbf{L} 的任一特征值满足 $|\lambda| < 1$, 故 SOR 方法收敛(注意当 $0 < \omega < 2$ 时, 可以证明 $(\omega + 2)^2 + \omega^2 > 0$).

定理 10 设 $\mathbf{Ax} = \mathbf{b}$, 如果:

(1) \mathbf{A} 为严格对角占优矩阵(或 \mathbf{A} 为弱对角占优不可约矩阵);

$$(2) 0 < \omega < 1.$$

则解 $\mathbf{Ax} = \mathbf{b}$ 的 SOR 迭代法收敛.

下面讨论迭代法的收敛速度.

由定理 3 证明中可知, 如果 $(\mathbf{B}) < 1$ 且 (\mathbf{B}) 越小时, 迭代法收敛越快. 现设有方程组

$$\mathbf{x} = \mathbf{Bx} + \mathbf{f}, \quad \mathbf{B} \in \mathbb{R}^{n \times n}$$

及一阶定常迭代法

$$\mathbf{x}^{(k+1)} = \mathbf{Bx}^{(k)} + \mathbf{f} \quad (k = 0, 1, \dots), \quad (3.9)$$

且设迭代法收敛, 记 $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^*$, 则 $\mathbf{x}^* = \mathbf{Bx}^* + \mathbf{f}$.

由基本定理有 $0 < (\mathbf{B}) < 1$, 且误差向量 $e^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}^*$ 满足

$$e^{(k)} = \mathbf{B}^{(k)} e^{(0)},$$

故

$$\mathbf{B}^{(k)} = \mathbf{B}^{(0)} \mathbf{B}^{(1)} \cdots \mathbf{B}^{(k-1)},$$

现设 \mathbf{B} 为对称矩阵, 则有

$$\begin{aligned} e^{(k)} &= \mathbf{B}^{(k)} e^{(0)} \\ &= [(\mathbf{B})]^k e^{(0)}. \end{aligned}$$

下面确定欲使初始误差缩小 10^{-s} 所需的迭代次数, 即欲使

$$[(\mathbf{B})]^k \leq 10^{-s}.$$

取对数, 得到所需最少迭代次数为

$$k = \frac{\ln 10}{\ln (\mathbf{B})}. \quad (3.10)$$

这说明, 所需迭代次数与 $R = -\ln (\mathbf{B})$ 成反比, 当 $(\mathbf{B}) < 1$ 越小,

$R = -\ln(\mathbf{B})$ 越大, 则 (3.10) 式满足所需迭代次数越少, 即迭代法收敛越快. 当迭代矩阵 \mathbf{B} 为一般矩阵时的讨论可参考 [7].

定义 5 称 $R(\mathbf{B}) = -\ln(\mathbf{B})$ 为迭代法 (3.9) 的渐近收敛速度, 简称迭代法收敛速度.

对于 SOR 迭代法希望选择松弛因子 ω 使迭代过程 (2.10) 收敛较快, 在理论上即确定 ω_{opt} 使

$$\min_{0 < \omega < 2} R(\mathbf{L}) = R(\mathbf{L}_{\text{opt}}).$$

对某些特殊类型的矩阵, 建立了 SOR 方法最佳松弛因子理论. 例如, 对所谓具有“性质 A”等条件的线性方程组建立了最佳松弛因子公式

$$\omega_{\text{opt}} = \frac{2}{1 + 1 - ((\mathbf{J}))^2},$$

其中 (\mathbf{J}) 为解 $\mathbf{Ax} = \mathbf{b}$ 的雅可比迭代法的迭代矩阵的谱半径.

在实际应用中, 对于某些椭圆型微分方程(模型问题)可以给出 ω_{opt} 的计算方法, 但一般来说, 计算 ω_{opt} 是有困难的, 可用试算的办法来确定一个适当的 ω .

算法 2(SOR 迭代法) 设 $\mathbf{Ax} = \mathbf{b}$, 其中 \mathbf{A} 为对称正定矩阵或为严格对角占优阵或为弱对角占优不可约矩阵等, 本算法用 SOR 迭代法求解 $\mathbf{Ax} = \mathbf{b}$, 数组 $x(n)$ 存放 $\mathbf{x}^{(0)}$ 及 $\mathbf{x}^{(k)}$, 用 $p_0 \max_{1 \leq i \leq n} |x_i| < \epsilon ps$ 控制迭代终止, 用 N_0 表示最大迭代次数.

1. $k = 0$

2. $x_i = 0.0 \quad (i = 1, 2, \dots, n)$

3. $k = k + 1$

4. $p_0 = 0.0$

5. 对于 $i = 1, 2, \dots, n$

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

(2) 如果 $|p| > p_0$ 则 $p_0 = |p|$

(3) $x_i = x_i + p$

6. 输出 p_0

7. 如果 $p_0 < \epsilon_{ps}$ 则输出 k , \mathbf{x} , 停机

8. 如果 $k < N_0$ 则转 3

9. 输出 N_0 及有关信息

[注] 可用 $\|\mathbf{r}^{(k)}\| < \epsilon_{ps}$ 来控制迭代终止, 其中 $\mathbf{r}^{(k)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}$.

6.4 分块迭代法

上述迭代法, 从 $\mathbf{x}^{(k)} = \mathbf{x}^{(k+1)}$ 计算时, 是逐个计算 $\mathbf{x}^{(k+1)}$ 的分量 $x_i^{(k+1)}$ ($i = 1, 2, \dots, n$), 这种迭代法又称为点迭代法, 下面介绍分块迭代法, 就是一块或一组未知数同时被改进.

设 $\mathbf{Ax} = \mathbf{b}$, 其中 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 为大型稀疏矩阵且将 \mathbf{A} 分块为三部分 $\mathbf{A} = \mathbf{D} - \mathbf{L} - \mathbf{U}$, 其中

$$\mathbf{A} = \begin{matrix} \mathbf{A}_1 & \mathbf{A}_2 & \dots & \mathbf{A}_q & \mathbf{A}_{11} \\ \mathbf{A}_{11} & \mathbf{A}_{12} & \dots & \mathbf{A}_{1q} & \mathbf{A}_{22} \\ \dots & \dots & \dots & \dots & \dots \\ \mathbf{A}_{q1} & \mathbf{A}_{q2} & \dots & \mathbf{A}_{qq} & \mathbf{A}_{qq} \\ \mathbf{0} & & & & \mathbf{0} - \mathbf{A}_2 - \dots - \mathbf{A}_q \\ -\mathbf{A}_{11} & \mathbf{0} & & & \mathbf{0} - \dots - \mathbf{A}_q \\ \dots & \dots & \dots & \dots & \dots \\ -\mathbf{A}_{q1} & -\mathbf{A}_{q2} & \dots & \mathbf{0} & \mathbf{0} \end{matrix}, \quad \mathbf{D} = \begin{matrix} \mathbf{A}_{11} \\ \mathbf{A}_{22} \\ \vdots \\ \mathbf{A}_{qq} \end{matrix}, \quad \mathbf{L} = \begin{matrix} \mathbf{0} \\ -\mathbf{A}_{11} \\ \dots \\ -\mathbf{A}_{q1} \end{matrix}, \quad \mathbf{U} = \begin{matrix} \mathbf{0} - \mathbf{A}_2 - \dots - \mathbf{A}_q \\ \dots \\ \mathbf{0} - \dots - \mathbf{A}_q \\ \dots \\ \mathbf{0} \end{matrix}.$$

且 \mathbf{A}_i ($i = 1, 2, \dots, q$) 为 $n_i \times n_i$ 非奇异矩阵, $n = \sum_{i=1}^q n_i$. 对 \mathbf{x} 及 \mathbf{b} 同样分块

$$\begin{array}{cc} \mathbf{x} & \mathbf{b} \\ \mathbf{x} = & , \quad \mathbf{b} = \\ \cdots & \cdots \\ \mathbf{x}_q & \mathbf{b}_q \end{array}$$

其中, $\mathbf{x} \in \mathbb{R}^{n_i}$, $\mathbf{b} \in \mathbb{R}^{n_i}$.

(1) 块雅可比迭代法(BJ)

选取分裂阵 \mathbf{M} 为 \mathbf{A} 的对角块部分, 即选

$$\mathbf{M} = \mathbf{D} \text{(块对角阵)},$$

$$\mathbf{A} = \mathbf{M} - \mathbf{N}.$$

于是, 得到块雅可比迭代法

$$\mathbf{x}^{(k+1)} = \mathbf{Bx}^{(k)} + \mathbf{f}, \quad (4.1)$$

其中迭代矩阵 $\mathbf{B} = \mathbf{I} - \mathbf{D}^{-1} \mathbf{A} = \mathbf{D}^{-1} (\mathbf{L} + \mathbf{U}) - \mathbf{J}$, $\mathbf{f} = \mathbf{D}^{-1} \mathbf{b}$.

或

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

由分块矩阵乘法, 得到块雅可比迭代法的具体形式

$$\mathbf{A}_{ii} \mathbf{x}_i^{(k+1)} = \mathbf{b} - \sum_{\substack{j=1 \\ j \neq i}}^q \mathbf{A}_{ij} \mathbf{x}_j^{(k)} \quad (i = 1, 2, \dots, q), \quad (4.2)$$

其中

$$\begin{array}{cc} \mathbf{x}^{(k)} \\ \mathbf{x}^{(k)} = & , \quad \mathbf{x}^{(k)} \in \mathbb{R}^{n_i} \\ \cdots \\ \mathbf{x}_q^{(k)} \end{array}$$

这说明, 块雅可比迭代法, 每迭代一步, 从 $\mathbf{x}^{(k)}$ 到 $\mathbf{x}^{(k+1)}$, 需要求解 q 个低阶方程组

$$\mathbf{A}_i \mathbf{x}_i^{(k+1)} = \mathbf{g}$$

$(i = 1, 2, \dots, q)$, 其中 \mathbf{g} 为(3.12)式右边部分.

(2) 块 SOR 迭代法(BSOR)

选取分裂矩阵 \mathbf{M} 为带松弛因子的 \mathbf{A} 块下三角部分, 即

$$\mathbf{M} = \frac{1}{\omega} (\mathbf{D} - \mathbf{L}),$$

$$\mathbf{A} = \mathbf{M} + \mathbf{N}.$$

得到块 SOR 迭代法

$$\mathbf{x}^{(k+1)} = \mathbf{L} \mathbf{x}^{(k)} + \mathbf{f}, \quad (4.3)$$

其中迭代矩阵

$$\begin{aligned}\mathbf{L} &= \mathbf{I} - (\mathbf{D} - \mathbf{L})^{-1} \mathbf{A} \\ &= (\mathbf{D} - \mathbf{L})^{-1} ((1 - \omega) \mathbf{D} + \mathbf{U}), \\ \mathbf{f} &= (\mathbf{D} - \mathbf{L})^{-1} \mathbf{b}.\end{aligned}$$

由分块矩阵乘法得到块 SOR 迭代法的具体形式

$$\begin{aligned}\mathbf{A}_{ii} \mathbf{x}_i^{(k+1)} &= \mathbf{A}_{ii} \mathbf{x}_i^{(k)} + \mathbf{b}_i - \sum_{j=1}^{i-1} \mathbf{A}_{ij} \mathbf{x}_j^{(k+1)} - \sum_{j=i}^q \mathbf{A}_{ij} \mathbf{x}_j^{(k)} \\ (i &= 1, 2, \dots, q; k = 0, 1, \dots),\end{aligned}$$

为松弛因子 .

(4.4)

于是, 当 $\mathbf{x}^{(k)}$ 及 $\mathbf{x}_j^{(k+1)}$ ($j = 1, 2, \dots, i - 1$) 已计算时, 解低阶方程组 (3.14) 可计算小块 $\mathbf{x}_i^{(k+1)}$. 从 $\mathbf{x}^{(k)}$ 到 $\mathbf{x}^{(k+1)}$ 共需要解 q 个低阶方程组, 当 \mathbf{A}_{ii} 为三对角阵或带状矩阵时, 可用直接法求解 .

我们给出下述结果 .

定理 11 设 $\mathbf{Ax} = \mathbf{b}$, 其中 $\mathbf{A} = \mathbf{D} - \mathbf{L} - \mathbf{U}$ (分块形式) .

(1) 如果 \mathbf{A} 为对称正定矩阵,

(2) $0 < \omega < 2$.

则解 $\mathbf{Ax} = \mathbf{b}$ 的 BSOR 迭代法收敛 .

评 注

本章介绍了解大型稀疏线性方程组的一些基本迭代法,例如,雅可比迭代法,高斯-塞德尔迭代法,及 SOR 迭代法,分块迭代法等. 并且建立了迭代法的一些基本理论. 在应用中 SOR 方法较为重要,它是解大型、稀疏线性方程组的有效方法之一.

迭代法是一种逐次逼近方法,在使用迭代法解方程组时,其系数矩阵在计算过程中始终不变.

迭代法具有循环的计算公式,方法简单,适宜解大型稀疏矩阵方程组,在计算机实现时只需存储 \mathbf{A} 的非零元素(或可按一定公式形成系数,这样 \mathbf{A} 就不需存储).

在使用迭代法时,要注意收敛性及收敛速度问题,使用 SOR 方法要选择较佳松弛因子.

迭代法的进一步学习,读者可参考文献[3],[14],[15],[16].

习 题

1. 设方程组

$$\begin{aligned} 5x_1 + 2x_2 + x_3 &= -12, \\ -x_1 + 4x_2 + 2x_3 &= 20, \\ 2x_1 - 3x_2 + 10x_3 &= 3. \end{aligned}$$

- (a) 考察用雅可比迭代法,高斯-塞德尔迭代法解此方程组的收敛性;
- (b) 用雅可比迭代法及高斯-塞德尔迭代法解此方程组,要求当 $\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| < 10^{-4}$ 时迭代终止.

2. 设方程组

$$\begin{array}{ll} (a) \quad 0.4x_1 + x_2 + 0.4x_3 = 1, & x_1 + 2x_2 - 2x_3 = 1, \\ 0.4x_1 + 0.8x_2 + x_3 = 2, & (b) \quad x_1 + x_2 + x_3 = 1, \\ 0.4x_1 + 0.8x_2 + x_3 = 3. & 2x_1 + 2x_2 + x_3 = 1. \end{array}$$

试考察解此方程组的雅可比迭代法及高斯-塞德尔迭代法的收敛性 .

3. 求证 $\lim_k \mathbf{A}_k = \mathbf{A}$ 的充要条件是对任何向量 \mathbf{x} 都有

$$\lim_k \mathbf{A}_k \mathbf{x} = \mathbf{A}\mathbf{x}.$$

4. 设 $\mathbf{Ax} = \mathbf{b}$, 其中 \mathbf{A} 对称正定, 问解此方程组的雅可比迭代法是否一定收敛? 试考察习题 2(a) 方程组 .

5. 用 SOR 方法解方程组(分别取松弛因子 $\omega = 1.03, \omega = 1, \omega = 1.1$)

$$\begin{aligned} 4x_1 - x_2 &= 1, \\ -x_1 + 4x_2 - x_3 &= 4, \\ -x_2 + 4x_3 &= -3. \end{aligned}$$

精确解 $\mathbf{x}^* = \begin{pmatrix} \frac{1}{2} \\ 1 \\ -\frac{1}{2} \end{pmatrix}^T$. 要求当 $\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| < 5 \times 10^{-6}$ 时迭代终止, 并且对每一个 ω 值确定迭代次数 .

6. 用 SOR 方法解方程组(取 $\omega = 0.9$)

$$\begin{aligned} 5x_1 + 2x_2 + x_3 &= -12, \\ -x_1 + 4x_2 + 2x_3 &= 20, \\ 2x_1 - 3x_2 + 10x_3 &= 3. \end{aligned}$$

要求当 $\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| < 10^{-4}$ 时迭代终止 .

7. 设有方程组 $\mathbf{Ax} = \mathbf{b}$, 其中 \mathbf{A} 为对称正定阵, 迭代公式

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + (\mathbf{b} - \mathbf{Ax}^{(k)}) \quad (k = 0, 1, 2, \dots),$$

试证明当 $0 < \rho < \frac{2}{\lambda_{\min}(\mathbf{A}) - \lambda_{\max}(\mathbf{A})}$ 时上述迭代法收敛(其中 $0 < \rho < \frac{2}{\lambda_{\min}(\mathbf{A}) - \lambda_{\max}(\mathbf{A})}$).

8. 证明矩阵

$$\mathbf{A} = \begin{matrix} 1 & a & a \\ a & 1 & a \\ a & a & 1 \end{matrix}$$

对于 $-\frac{1}{2} < a < 1$ 是正定的, 而雅可比迭代只对 $-\frac{1}{2} < a < \frac{1}{2}$ 是收敛的 .

9. 给定迭代过程 $\mathbf{x}^{(k+1)} = \mathbf{G}\mathbf{x}^{(k)} + \mathbf{g}$, 其中 $\mathbf{G} \in \mathbb{R}^{n \times n}$ ($k = 0, 1, 2, \dots$), 试证明: 如果 \mathbf{G} 的特征值 $\lambda_i(\mathbf{G}) = 0$ ($i = 1, 2, \dots, n$), 则此迭代过程最多迭代 n 次收敛于方程组的解 .

10. 设 \mathbf{A} 为严格对角占优阵, 且 $0 < \omega < 1$. 求证解 $\mathbf{Ax} = \mathbf{b}$ 的 SOR 迭代法收敛 .

第 7 章 非线性方程求根

7.1 方程求根与二分法

7.1.1 引言

本章主要讨论单变量非线性方程

$$f(x) = 0 \quad (1.1)$$

的求根问题, 这里 $x \in \mathbf{R}$, $f(x) \in C[a, b]$. 在科学与工程计算中有大量方程求根问题, 其中一类特殊的问题是多项式方程

$$f(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0 \quad (a_n \neq 0), \quad (1.2)$$

其中系数 a_i ($i=0, 1, \dots, n$) 为实数.

方程 $f(x) = 0$ 的根 x^* , 又称为函数 $f(x)$ 的零点, 它使 $f(x^*) = 0$, 若 $f(x)$ 可分解为

$$f(x) = (x - x^*)^m g(x),$$

其中 m 为正整数, 且 $g(x^*) \neq 0$. 当 $m=1$ 时, 则称 x^* 为单根, 若 $m > 1$ 称 x^* 为(1.1)的 m 重根, 或 x^* 为 $f(x)$ 的 m 重零点. 若 x^* 是 $f(x)$ 的 m 重零点, 且 $g(x)$ 充分光滑, 则

$$f(x^*) = f'(x^*) = \dots = f^{(m-1)}(x^*) = 0, \quad f^{(m)}(x^*) \neq 0.$$

当 $f(x)$ 为代数多项式(1.2)时, 根据代数基本定理可知, n 次方程在复数域有且只有 n 个根(含复根, m 重根为 m 个根), $n=1, 2$ 时方程的根是大家熟悉的, $n=3, 4$ 时虽有求根公式但比较复杂, 可在数学手册中查到, 但已不适合于数值计算, 而 $n \geq 5$ 时就不能用公式表示方程的根. 因此, 通常对 $n \geq 3$ 的多项式方程求根与一般连续函数方程(1.1)一样都可采用迭代法求根. 迭代法要求先

给出根 x^* 的一个近似, 若 $f(x) \in C[a, b]$ 且 $f(a)f(b) < 0$, 根据连续函数性质可知 $f(x) = 0$ 在 (a, b) 内至少有一个实根, 这时称 $[a, b]$ 为方程(1.1)的有根区间. 通常可通过逐次搜索法求得方程(1.1)的有根区间.

例 1 求方程 $f(x) = x^3 - 11.1x^2 + 38.8x - 41.77 = 0$ 的有根区间.

解 根据有根区间定义, 对 $f(x) = 0$ 的根进行搜索计算, 结果如下:

x	0	1	2	3	4	5	6
$f(x)$ 的符号	-	-	+	+	-	-	+

由此可知方程的有根区间为 $[1, 2], [3, 4], [5, 6]$.

7.1.2 二分法

考察有根区间 $[a, b]$, 取中点 $x_0 = (a + b)/2$ 将它分为两半, 假设中点 x_0 不是 $f(x)$ 的零点, 然后进行根的搜索, 即检查 $f(x_0)$ 与 $f(a)$ 是否同号, 如果确系同号, 说明所求的根 x^* 在 x_0 的右侧, 这时令 $a = x_0, b = b$; 否则 x^* 必在 x_0 的左侧, 这时令 $a = a, b = x_0$ (图 7-1). 不管出现哪一种情况, 新的有根区间 $[a, b]$ 的长度仅为 $[a, b]$ 的一半.

对压缩了的有根区间 $[a, b]$ 又可施行同样的手续, 即用中点 $x_1 = (a_1 + b_1)/2$ 将区间 $[a_1, b_1]$ 再分为两半, 然后通过根的搜索判定所求的根在 x_1 的哪一侧, 从而又确定一个新的有根区间 $[a_1, b_1]$, 其长度是 $[a, b]$ 的一半.

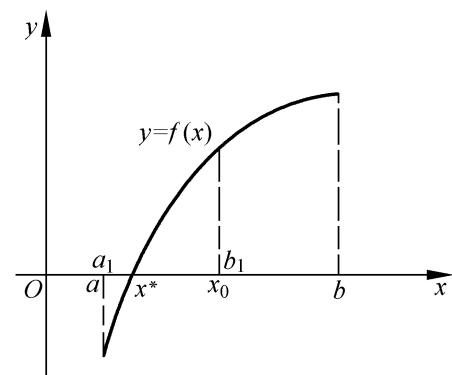


图 7-1

如此反复二分下去, 即可得出一系列有根区间

$$[a, b] \quad [a_1, b_1] \quad [a_2, b_2] \quad \dots \quad [a_k, b_k] \quad \dots,$$

其中每个区间都是前一个区间的一半, 因此 $[a_k, b_k]$ 的长度

$$b_k - a_k = (b - a)/2^k$$

当 k 时趋于零, 就是说, 如果二分过程无限地继续下去, 这些区间最终必收缩于一点 x^* , 该点显然就是所求的根.

每次二分后, 设取有根区间 $[a_k, b_k]$ 的中点

$$x_k = (a_k + b_k)/2$$

作为根的近似值, 则在二分过程中可以获得一个近似根的序列

$$x_0, x_1, x_2, \dots, x_k, \dots,$$

该序列必以根 x^* 为极限.

不过在实际计算时, 我们不可能完成这个无限过程, 其实也没有这种必要, 因为数值分析的结果允许带有一定的误差. 由于

$$|x^* - x_k| = (b_k - a_k)/2 = (b - a)/2^{k+1}, \quad (1.3)$$

只要二分足够多次(即 k 充分大), 便有

$$|x^* - x_k| < ,$$

这里 为预定的精度.

例 2 求方程

$$f(x) = x^3 - x - 1 = 0$$

在区间 $[1.0, 1.5]$ 内的一个实根, 要求准确到小数点后的第 2 位.

解 这里 $a = 1.0$, $b = 1.5$, 而 $f(a) < 0$, $f(b) > 0$. 取 $[a, b]$ 的中点 $x_0 = 1.25$, 将区间二等分, 由于 $f(x_0) < 0$, 即 $f(x_0)$ 与 $f(a)$ 同号, 故所求的根 x^* 必在 x_0 右侧, 这时应令 $a = x_0 = 1.25$, $b = b = 1.5$, 而得到新的有根区间 $[a_1, b]$.

如此反复二分下去. 二分过程无需赘述. 我们现在预估所要二分的次数, 按误差估计(1.3)式, 只要二分 6 次($k=6$), 便能达到预定的精度

$$|x^* - x_6| < 0.005 .$$

二分法的计算结果如表 7-1.

二分法是计算机上的一种常用算法,下面列出计算步骤:

表 7-1 计算结果

k	a_k	b_k	x_k	$f(x_k)$ 符号
0	1.0	1.5	1.25	-
1	1.25		1.375	+
2		1.375	1.3125	-
3	1.3125		1.3438	+
4		1.3438	1.3281	+
5		1.3281	1.3203	-
6	1.3203		1.3242	-

步骤 1 准备 计算 $f(x)$ 在有根区间 $[a, b]$ 端点处的值 $f(a), f(b)$.

步骤 2 二分 计算 $f(x)$ 在区间中点 $\frac{a+b}{2}$ 处的值 $f\left(\frac{a+b}{2}\right)$.

步骤 3 判断 若 $f\left(\frac{a+b}{2}\right) = 0$, 则 $\frac{a+b}{2}$ 即是根, 计算过程结束, 否则检验.

若 $f\left(\frac{a+b}{2}\right) f(a) < 0$, 则以 $\frac{a+b}{2}$ 代替 b , 否则以 $\frac{a+b}{2}$ 代替 a .

反复执行步骤 2 和步骤 3, 直到区间 $[a, b]$ 长度小于允许误差, 此时中点 $\frac{a+b}{2}$ 即为所求近似根.

上述二分法的优点是算法简单, 且总是收敛的, 缺点是收敛太慢, 故一般不单独将其用于求根, 只用其为根求得一个较好的近似值.

7.2 迭代法及其收敛性

7.2.1 不动点迭代法

将方程(1.1)改写成等价的形式

$$x = \varphi(x). \quad (2.1)$$

若要求 x^* 满足 $f(x^*) = 0$, 则 $x^* = \varphi(x^*)$; 反之亦然, 称 x^* 为函数 $\varphi(x)$ 的一个不动点. 求 $f(x)$ 的零点就等价于求 $\varphi(x)$ 的不动点, 选择一个初始近似值 x_0 , 将它代入(2.1)右端, 即可求得

$$x_1 = \varphi(x_0).$$

可以如此反复迭代计算

$$x_{k+1} = \varphi(x_k) \quad (k = 0, 1, \dots). \quad (2.2)$$

$\varphi(x)$ 称为迭代函数. 如果对任何 $x_0 \in [a, b]$, 由(2.2)得到的序列 $\{x_k\}$ 有极限

$$\lim_k x_k = x^*.$$

则称迭代方程(2.2)收敛, 且 $x^* = \varphi(x^*)$ 为 $\varphi(x)$ 的不动点, 故称(2.2)为不动点迭代法.

上述迭代法是一种逐次逼近法, 其基本思想是将隐式方程(2.1)归结为一组显式的计算公式(2.2), 就是说, 迭代过程实质上是一个逐步显示化的过程.

我们用几何图像来显示迭代过程. 方程 $x = \varphi(x)$ 的求根问题在 xy 平面上就是要确定曲线 $y = \varphi(x)$ 与直线 $y = x$ 的交点 P^* (图 7-2). 对于 x^*

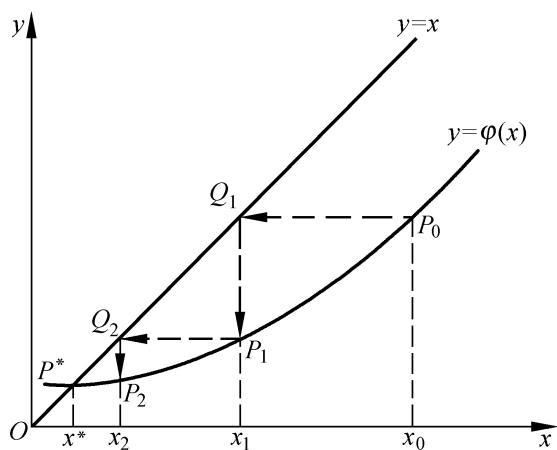


图 7-2

的某个近似值 x_0 , 在曲线 $y = f(x)$ 上可确定一点 P_0 , 它以 x_0 为横坐标, 而纵坐标则等于 $f(x_0) = x_1$. 过 P_0 引平行于 x 轴的直线, 设此直线交直线 $y = x$ 于点 Q_1 , 然后过 Q_1 再作平行于 y 轴的直线, 它与曲线 $y = f(x)$ 的交点记作 P_1 , 则点 P_1 的横坐标为 x_1 , 纵坐标则等于 $f(x_1) = x_2$. 按图 7-2 中箭头所示的路径继续做下去, 在曲线 $y = f(x)$ 上得到点列 P_1, P_2, \dots , 其横坐标分别为依公式 $x_{k+1} = f(x_k)$ 求得的迭代值 x_1, x_2, \dots . 如果点列 $\{P_k\}$ 趋向于点 P^* , 则相应的迭代值 x_k 收敛到所求的根 x^* .

例 3 求方程

$$f(x) = x^3 - x - 1 = 0 \quad (2.3)$$

在 $x_0 = 1.5$ 附近的根 x^* .

解 设将方程(2.3)改写成下列形式

$$x = \sqrt[3]{x + 1}.$$

据此建立迭代公式

$$x_{k+1} = \sqrt[3]{x_k + 1} \quad (k = 0, 1, 2, \dots).$$

表 7-2 记录了各步迭代的结果.

我们看到, 如果仅取 6 位数字, 那么结果 x_7 与 x_8 完全相同, 这时可以认为 x_7 实际上已满足方程(2.3), 即为所求的根.

应当指出, 迭代法的效果并不是总能令人满意的. 譬如, 设依方程(2.3)的另一种等价形式

$$x = x^3 - 1$$

建立迭代公式

$$x_{k+1} = \sqrt[3]{x_k^3 - 1}.$$

迭代初值仍取 $x_0 = 1.5$, 则有

表 7-2 计算结果

k	x_k	k	x_k
0	1.5	5	1.32476
1	1.35721	6	1.32473
2	1.33086	7	1.32472
3	1.32588	8	1.32472
4	1.32494		

$$x_1 = 2.375, \quad x_2 = 12.39.$$

继续迭代下去已经没有必要, 因为结果显然会越来越大, 不可能趋于某个极限. 这种不收敛的迭代过程称作是发散的. 一个发散的迭代过程, 纵使进行了千百次迭代, 其结果也是毫无价值的.

例 3 表明原方程化为(2.1)的形式不同, 有的收敛, 有的发散, 只有收敛的迭代过程(2.2)才有意义, 为此我们首先要研究 (x) 的不动点的存在性及迭代法(2.2)的收敛性.

7.2.2 不动点的存在性与迭代法的收敛性

首先考察 (x) 在 $[a, b]$ 上不动点的存在唯一性.

定理 1 设 $(x) \in C[a, b]$ 满足以下两个条件:

1° 对任意 $x \in [a, b]$ 有 $a \leq (x) \leq b$.

2° 存在正常 $L < 1$, 使对任意 $x, y \in [a, b]$ 都有

$$|(x) - (y)| \leq L |x - y|. \quad (2.4)$$

则 (x) 在 $[a, b]$ 上存在唯一的不动点 x^* .

证明 先证不动点存在性. 若 $(a) = a$ 或 $(b) = b$, 显然 (x) 在 $[a, b]$ 上存在不动点. 因 $a \leq (x) \leq b$, 以下设 $(a) > a$ 及 $(b) < b$, 定义函数

$$f(x) = (x) - x.$$

显然 $f(x) \in C[a, b]$, 且满足 $f(a) = (a) - a > 0$, $f(b) = (b) - b < 0$, 由连续函数性质可知存在 $x^* \in (a, b)$ 使 $f(x^*) = 0$, 即 $x^* = (x^*)$, x^* 即为 (x) 的不动点.

再证唯一性. 设 x_1^* 及 $x_2^* \in [a, b]$ 都是 (x) 的不动点, 则由(2.4)得

$$\begin{aligned} |x_1^* - x_2^*| &= |(x_1^*) - (x_2^*)| \\ &\leq L |x_1^* - x_2^*| < |x_1^* - x_2^*|. \end{aligned}$$

引出矛盾. 故 (x) 的不动点只能是唯一的.

证毕.

在 (x) 的不动点存在唯一的情况下, 可得到迭代法(2.2)收敛的一个充分条件.

定理2 设 $(x) \in C[a, b]$ 满足定理1中的两个条件, 则对任意 $x_0 \in [a, b]$, 由(2.2)得到的迭代序列 $\{x_k\}$ 收敛到 (x) 的不动点 x^* , 并有误差估计

$$|x_k - x^*| \leq \frac{L^k}{1-L} |x_1 - x_0|. \quad (2.5)$$

证明 设 $x^* \in [a, b]$ 是 (x) 在 $[a, b]$ 上的唯一不动点, 由条件1°, 可知 $\{x_k\} \subset [a, b]$, 再由(2.4)得

$$\begin{aligned} |x_k - x^*| &= |(x_{k-1}) - (x^*)| \\ &\leq L |x_{k-1} - x^*| \dots L^k |x_0 - x^*|. \end{aligned}$$

因 $0 < L < 1$, 故当 k 时序列 $\{x_k\}$ 收敛到 x^* .

下面再证明估计式(2.5), 由(2.4)有

$$|x_{k+1} - x_k| = |(x_k) - (x_{k-1})| \leq L |x_k - x_{k-1}|. \quad (2.6)$$

据此反复递推得

$$|x_{k+1} - x_k| \leq L^k |x_1 - x_0|.$$

于是对任意正整数 p 有

$$\begin{aligned} |x_{k+p} - x_k| &= |x_{k+p} - x_{k+p-1}| + |x_{k+p-1} - x_{k+p-2}| \\ &\quad + \dots + |x_{k+1} - x_k| \\ &\leq (L^{k+p-1} + L^{k+p-2} + \dots + L^k) |x_1 - x_0| \\ &\leq \frac{L^k}{1-L} |x_1 - x_0|. \end{aligned}$$

在上式令 $p \rightarrow \infty$, 注意到 $\lim_{p \rightarrow \infty} x_{k+p} = x^*$ 即得式(2.5). 证毕.

迭代过程是个极限过程. 在用迭代法进行实际计算时, 必须按精度要求控制迭代次数. 误差估计式(2.5)原则上可用于确定迭代次数, 但它由于含有信息 L 而不便于实际应用. 根据式(2.6), 对任意正整数 p 有

$$|x_{k+p} - x_k| \leq (L^{p-1} + L^{p-2} + \dots + 1) |x_{k+1} - x_k|$$

$$\frac{1}{1 - L} / |x_{k+1} - x_k|.$$

在上式中令 p 知

$$|x^* - x_k| / \frac{1}{1 - L} / |x_{k+1} - x_k|.$$

由此可见, 只要相邻两次计算结果的偏差 $|x_{k+1} - x_k|$ 足够小即可保证近似值 x_k 具有足够精度 .

对定理 1 和定理 2 中的条件 2°, 在使用时如果 $(x) \in C[a, b]$ 且对任意 $x \in [a, b]$ 有

$$|(x)| / L < 1, \quad (2.7)$$

则由中值定理可知对 " $x, y \in [a, b]$ " 有

$$|(x) - (y)| = |()(x - y)| / L / |x - y|, \quad (a, b).$$

它表明定理中的条件 2° 可用 (2.7) 代替 .

在例 3 中, 当 $(x) = x^3 + 1$ 时, $(x) = \frac{1}{3}(x+1)^{2/3}$, 在区间 $[1, 2]$ 中, $|(x)| = \frac{1}{3} \frac{1}{4}^{1/3} < 1$, 故 (2.7) 成立 . 又因 $1^3 = 2$, $(x) = 3^3 = 2$, 故定理 1 中条件 1° 也成立 . 所以迭代法是收敛的 . 而当 $(x) = x^3 - 1$ 时, $(x) = 3x^2$ 在区间 $[1, 2]$ 中 $|(x)| > 1$ 不满足定理条件 .

7.2.3 局部收敛性与收敛阶

上面给出了迭代序列 $\{x_k\}$ 在区间 $[a, b]$ 上的收敛性, 通常称为全局收敛性 . 有时不易检验定理的条件, 实际应用时通常只在不动点 x^* 的邻近考察其收敛性, 即局部收敛性 .

定义 1 设 (x) 有不动点 x^* , 如果存在 x^* 的某个邻域 $R: |x - x^*| < R$, 对任意 $x_0 \in R$, 迭代 (2.2) 产生的序列 $\{x_k\} \subset R$, 且收敛到 x^* , 则称迭代法 (2.2) 局部收敛 .

定理 3 设 x^* 为 (x) 的不动点, (x) 在 x^* 的某个邻域连

续,且 $|f(x^*)| < 1$,则迭代法(2.2)局部收敛.

证明 由连续函数的性质,存在 x^* 的某个邻域 $R: |x - x^*|$,使对于任意 $x \in R$ 成立

$$|f'(x)| < L < 1.$$

此外,对于任意 $x \in R$,总有 $|f(x)| < R$,这是因为

$$|f(x) - f(x^*)| = |f(x) - f(x^*)| / L \cdot |x - x^*| / |x - x^*|.$$

于是依据定理2可以断定迭代过程 $x_{k+1} = f(x_k)$ 对于任意初值 $x_0 \in R$ 均收敛.证毕.

下面讨论迭代序列的收敛速度问题,先看例.

例4 用不同方法求方程 $x^2 - 3 = 0$ 的根 $x^* = \sqrt{3}$.

解 这里 $f(x) = x^2 - 3$,可改写为各种不同的等价形式 $x = f(x)$,其不动点为 $x^* = \sqrt{3}$.由此构造不同的迭代法:

$$(1) \quad x_{k+1} = x_k^2 + x_k - 3, \quad f(x) = x^2 + x - 3, \quad f'(x) = 2x + 1,$$

$$(x^*) = f(\sqrt{3}) = 2\sqrt{3} + 1 > 1.$$

$$(2) \quad x_{k+1} = \frac{3}{x_k}, \quad f(x) = \frac{3}{x}, \quad f'(x) = -\frac{3}{x^2}, \quad (x^*) = -1.$$

$$(3) \quad x_{k+1} = x_k - \frac{1}{4}(x_k^2 - 3), \quad f(x) = x - \frac{1}{4}(x^2 - 3),$$

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

$$(x^*) = 1 - \frac{3}{2} = 0.134 < 1.$$

$$(4) \quad x_{k+1} = \frac{1}{2}x_k + \frac{3}{x_k}, \quad f(x) = \frac{1}{2}x + \frac{3}{x},$$

$$f'(x) = \frac{1}{2} - \frac{3}{x^2}, \quad (x^*) = f(\sqrt{3}) = 0.$$

取 $x_0 = 2$,对上述4种迭代法,计算三步所得的结果如下表.

表 7-3 计算结果

k	x_k	迭代法(1)	迭代法(2)	迭代法(3)	迭代法(4)
0	x_0	2	2	2	2
1	x_1	3	1.5	1.75	1.75
2	x_2	9	2	1.73475	1.732143
3	x_3	87	1.5	1.732361	1.732051
...

注意 $3 = 1.7320508\dots$, 从计算结果看到迭代法(1)及(2)均不收敛, 且它们均不满足定理 3 中的局部收敛条件, 迭代法(3)和(4)均满足局部收敛条件, 且迭代法(4)比(3)收敛快, 因在迭代法(4)中 $(x^*) = 0$. 为了衡量迭代法(2.2)收敛速度的快慢可给出以下定义.

定义 2 设迭代过程 $x_{k+1} = (x_k)$ 收敛于方程 $x = (x)$ 的根 x^* , 如果迭代误差 $e_k = x_k - x^*$ 当 k 时成立下列渐近关系式

$$\frac{e_{k+1}}{e_k^p} \quad C \quad (\text{常数 } C > 0),$$

则称该迭代过程是 p 阶收敛的. 特别地, $p = 1$ 时称线性收敛, $p > 1$ 时称超线性收敛, $p = 2$ 时称平方收敛.

定理 4 对于迭代过程 $x_{k+1} = (x_k)$, 如果 (x) 在所求根 x^* 的邻近连续, 并且

$$(x^*) = (x^*) = \dots = {}^{(p-1)}(x^*) = 0, \\ {}^{(p)}(x^*) \neq 0, \quad (2.8)$$

则该迭代过程在点 x^* 邻近是 p 阶收敛的.

证明 由于 $(x^*) = 0$, 据定理 3 立即可以断定迭代过程 $x_{k+1} = (x_k)$ 具有局部收敛性.

再将 (x_k) 在根 x^* 处做泰勒展开, 利用条件(2.8), 则有

$$(x_k) = (x^*) + \frac{^{(p)}(\cdot)}{p!} (x_k - x^*)^p, \text{ 在 } x_k \text{ 与 } x^* \text{ 之间.}$$

注意到 $(x_k) = x_{k+1}$, $(x^*) = x^*$, 由上式得

$$x_{k+1} - x^* = \frac{^{(p)}(\cdot)}{p!} (x_k - x^*)^p,$$

因此对迭代误差, 当 k 时有

$$\frac{e_{k+1}}{e_k^p} = \frac{^{(p)}(x^*)}{p!}. \quad (2.9)$$

这表明迭代过程 $x_{k+1} = (x_k)$ 确实为 p 阶收敛. 证毕.

上述定理告诉我们, 迭代过程的收敛速度依赖于迭代函数 (x) 的选取. 如果当 $x \in [a, b]$ 时 $(x) = 0$, 则该迭代过程只可能是线性收敛.

在例 4 中, 迭代法(3)的 $(x^*) = 0$, 故它只是线性收敛, 而迭代法(4)的 $(x^*) = 0$, 而 $(x) = \frac{6}{x^3}$, $(x^*) = \frac{2}{3} - 0$. 由定理 4 知 $p = 2$, 即该迭代过程为 2 阶收敛.

7.3 迭代收敛的加速方法

7.3.1 埃特金加速收敛方法

对于收敛的迭代过程, 只要迭代足够多次, 就可以使结果达到任意的精度, 但有时迭代过程收敛缓慢, 从而使计算量变得很大, 因此迭代过程的加速是个重要的课题.

设 x_0 是根 x^* 的某个近似值, 用迭代公式校正一次得

$$x_1 = (x_0),$$

而由微分中值定理, 有

$$x_1 - x^* = (x_0) - (x^*) = (\cdot)(x_0 - x^*),$$

其中 \cdot 介于 x^* 与 x_0 之间.

假定 (x) 改变不大, 近似地取某个近似值 L , 则有

$$x_1 - x^* = L(x_0 - x^*). \quad (3.1)$$

若将校正值 $x_1 = (x_0)$ 再校正一次, 又得

$$x_2 = (x_1),$$

由于

$$x_2 - x^* = L(x_1 - x^*),$$

将它与(3.1)式联立, 消去未知的 L , 有

$$\frac{x_1 - x^*}{x_2 - x^*} = \frac{x_0 - x^*}{x_1 - x^*}.$$

由此推知

$$x^* = \frac{x_0 x_2 - x_1^2}{x_2 - 2x_1 + x_0} = x_0 - \frac{(x_1 - x_0)^2}{x_2 - 2x_1 + x_0}.$$

在计算了 x_1 及 x_2 之后, 可用上式右端作为 x^* 的新近似, 记作 \bar{x} . 一般情形是由 x_k 计算 x_{k+1} , x_{k+2} , 记

$$\begin{aligned} \bar{x}_{k+1} &= x_k - \frac{(x_{k+1} - x_k)^2}{x_k - 2x_{k+1} + x_{k+2}} \\ &= x_k - (x_k)^2 / x_{k+1}^2 \quad (k = 0, 1, \dots). \end{aligned} \quad (3.2)$$

(3.2) 称为埃特金(Aitken)² 加速方法.

可以证明

$$\lim_k \frac{\bar{x}_{k+1} - x^*}{x_k - x^*} = 0.$$

它表明序列 $\{\bar{x}_k\}$ 的收敛速度比 $\{x_k\}$ 的收敛速度快.

7.3.2 斯蒂芬森迭代法

埃特金方法不管原序列 $\{x_k\}$ 是怎样产生的, 对 $\{x_k\}$ 进行加速计算, 得到序列 $\{\bar{x}_k\}$. 如果把埃特金加速技巧与不动点迭代结合, 则可得到如下的迭代法:

$$\begin{aligned} y_k &= (x_k), \quad z_k = (y_k), \\ x_{k+1} &= x_k - \frac{(y_k - x_k)^2}{z_k - 2y_k + x_k} \quad (k = 0, 1, \dots), \end{aligned} \quad (3.3)$$

称为斯蒂芬森(Steffensen)迭代法. 它可以这样理解, 我们要求 $x = (x)$ 的根 x^* , 令 $(x) = (x) - x$, $(x^*) = (x^*) - x^* = 0$, 已知 x^* 的近似值 x_k 及 y_k , 其误差分别为

$$(x_k) = (x_k) - x_k = y_k - x_k,$$

$$(y_k) = (y_k) - y_k = z_k - y_k.$$

把误差 (x) “外推到零”, 即过 $(x_k, (x_k))$ 及 $(y_k, (y_k))$ 两点做线性插值函数, 它与 x 轴交点就是(3.3)中的 x_{k+1} , 即方程

$$(x_k) + \frac{(y_k) - (x_k)}{y_k - x_k} (x - x_k) = 0$$

的解

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

实际上(3.3)是将不动点迭代法(2.2)计算两步合并成一步得到的, 可将它写成另一种不动点迭代

$$x_{k+1} = (x_k) \quad (k = 0, 1, \dots), \quad (3.4)$$

其中

$$(x) = x - \frac{[(x) - x]^2}{((x)) - 2(x) + x}. \quad (3.5)$$

对不动点迭代(3.4)有以下局部收敛性定理.

定理5 若 x^* 为(3.5)定义的迭代函数 (x) 的不动点, 则 x^* 为 (x) 的不动点. 反之, 若 x^* 为 (x) 的不动点, 设 (x) 存在, $(x^*) = 1$, 则 x^* 是 (x) 的不动点, 且斯蒂芬森迭代法(3.3)是2阶收敛的.

证明可见[2].

例5 用斯蒂芬森迭代法求解方程(2.3).

解 例3中已指出下列迭代

$$x_{k+1} = x_k^3 - 1$$

是发散的, 现用(3.3)计算, 取 $(x) = x^3 - 1$, 计算结果如下表.

表 7-4 计算结果

k	x_k	y_k	z_k
0	1.5	2.37500	12.3965
1	1.41629	1.84092	5.23888
2	1.35565	1.49140	2.31728
3	1.32895	1.34710	1.44435
4	1.32480	1.32518	1.32714
5	1.32472		

计算表明它是收敛的, 这说明即使迭代法(2.2)不收敛, 用斯蒂芬森迭代法(3.3)仍可能收敛. 至于原来已收敛的迭代法(2.2), 由定理 5 可知它可达到 2 阶收敛. 更进一步还可知若(2.2)为 p 阶收敛, 则(3.3)为 $p+1$ 阶收敛(见[2]).

例 6 求方程 $3x^2 - e^x = 0$ 在 $[3, 4]$ 中的解.

解 由方程得 $e^x = 3x^2$, 取对数得

$$x = \ln 3x^2 = 2\ln x + \ln 3 = (x).$$

若构造迭代法

$$x_{k+1} = 2\ln x_k + \ln 3,$$

由于 $(x) = \frac{2}{x}$, $\max_{3 \leq x \leq 4} |(x)| = \frac{2}{3} < 1$, 且当 $x \in [3, 4]$ 时, (x)

$[3, 4]$, 根据定理 2 此迭代法是收敛的. 若取 $x_0 = 3.5$ 迭代 16 次得 $x_{16} = 3.73307$, 有六位有效数字.

若用(3.3)进行加速, 计算结果如下:

k	x_k	y_k	z_k
0	3.5	3.60414	3.66202
1	3.73444	3.73381	3.73347
2	3.73307		

这里计算2步(相当于(2.2)迭代4步)结果与 x_{16} 相同,说明用迭代法(3.3)的收敛速度比迭代法(2.2)快得多.

7.4 牛顿法

7.4.1 牛顿法及其收敛性

对于方程 $f(x)=0$,如果 $f(x)$ 是线性函数,则它的求根是容易的.牛顿法实质上是一种线性化方法,其基本思想是将非线性方程 $f(x)=0$ 逐步归结为某种线性方程来求解.

设已知方程 $f(x)=0$ 有近似根 x_k (假定 $f'(x_k)\neq 0$),将函数 $f(x)$ 在点 x_k 展开,有

$$f(x) = f(x_k) + f'(x_k)(x - x_k),$$

于是方程 $f(x)=0$ 可近似地表示为

$$f(x_k) + f'(x_k)(x - x_k) = 0.$$

(4.1)

这是个线性方程,记其根为 x_{k+1} ,则 x_{k+1} 的计算公式为

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} \quad (k=0,1,\dots),$$

(4.2)

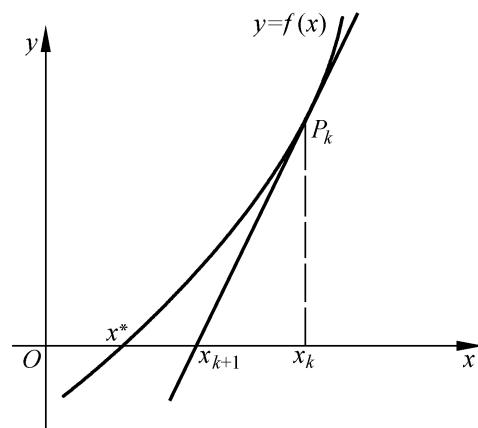


图 7-3

这就是牛顿(Newton)法.

牛顿法有明显的几何解释.方程 $f(x)=0$ 的根 x^* 可解释为曲线 $y=f(x)$ 与 x 轴的交点的横坐标(图7-3).设 x_k 是根 x^* 的某个近似值,过曲线 $y=f(x)$ 上横坐标为 x_k 的点 P_k 引切线,并将该切线与 x 轴的交点的横坐标 x_{k+1} 作为 x^* 的新的近似值.注意到切线方程为

$$y = f(x_k) + f'(x_k)(x - x_k).$$

这样求得的值 x_{k+1} 必满足(4.1), 从而就是牛顿公式(4.2)的计算结果. 由于这种几何背景, 牛顿法亦称切线法.

关于牛顿法(4.2)的收敛性, 可直接由定理4得到, 对(4.2)其迭代函数为

$$(x) = x - \frac{f(x)}{f'(x)},$$

由于

$$(x) = \frac{f(x)f''(x)}{[f'(x)]^2}.$$

假定 x^* 是 $f(x)$ 的一个单根, 即 $f(x^*)=0$, $f'(x^*) \neq 0$, 则由上式知 $(x^*)=0$, 于是依据定理4可以断定, 牛顿法在根 x^* 的邻近是平方收敛的. 又因

$$(x^*) = \frac{f(x^*)}{f'(x^*)}, \text{故由(2.9)可得}$$

$$\lim_k \frac{x_{k+1} - x^*}{(x_k - x^*)^2} = \frac{f(x^*)}{2f'(x^*)}. \quad (4.3)$$

例7 用牛顿法解方程

$$xe^x - 1 = 0. \quad (4.4)$$

解 这里牛顿公式为

$$x_{k+1} = x_k - \frac{x_k - e^{-x_k}}{1 + x_k},$$

取迭代初值 $x_0 = 0.5$, 迭代结果

列于表7-5中.

所给方程(4.4)实际上是方程 $x = e^{-x}$ 的等价形式. 若用不动点迭代到同一精度要迭代17次, 可见牛顿法的收敛速度是很快的.

下面列出牛顿法的计算

表7-5 计算结果

k	x_k
0	0.5
1	0.57102
2	0.56716
3	0.56714

步骤：

步骤 1 准备 选定初始近似值 x_0 , 计算 $f_0 = f(x_0)$,
 $f'_0 = f'(x_0)$.

步骤 2 迭代 按公式

$$x_1 = x_0 - \frac{f_0}{f'_0}$$

迭代一次, 得新的近似值 x_1 , 计算 $f_1 = f(x_1)$, $f'_1 = f'(x_1)$.

步骤 3 控制 如果 $|x_1 - x_0| < \epsilon_1$ 或 $|f_1| < \epsilon_2$, 则终止迭代, 以 x_1 作为所求的根; 否则转步骤 4. 此处 ϵ_1, ϵ_2 是允许误差, 而

$$|x_1 - x_0|, \quad \text{当 } |x_1| < C \text{ 时};$$

$$= \frac{|x_1 - x_0|}{|x_1|}, \quad \text{当 } |x_1| \geq C \text{ 时},$$

其中 C 是取绝对误差或相对误差的控制常数, 一般可取 $C = 1$.

步骤 4 修改 如果迭代次数达到预先指定的次数 N , 或者 $f_1 = 0$, 则方法失败; 否则以 (x_1, f_1, f'_1) 代替 (x_0, f_0, f'_0) 转步骤 2 继续迭代.

7.4.2 牛顿法应用举例

对于给定的正数 C , 应用牛顿法解二次方程

$$x^2 - C = 0,$$

可导出求开方值 C 的计算程序

$$x_{k+1} = \frac{1}{2} x_k + \frac{C}{x_k}. \quad (4.5)$$

我们现在证明, 这种迭代公式对于任意初值 $x_0 > 0$ 都是收敛的.

事实上, 对(4.5)式施行配方手续, 易知

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

$$x_{k+1} + C = \frac{1}{2x_k} (x_k + C)^2.$$

以上两式相除得

$$\frac{x_{k+1} - C}{x_{k+1} + C} = \frac{x_k - C}{x_k + C}^2.$$

据此反复递推有

$$\frac{x_k - C}{x_k + C} = \frac{x_0 - C}{x_0 + C}^{2^k}, \quad (4.6)$$

记

$$q = \frac{x_0 - C}{x_0 + C},$$

整理(4.6)式, 得

$$x_k - C = 2^k C \frac{q^{2^k}}{1 - q^{2^k}}.$$

对任意 $x_0 > 0$, 总有 $|q| < 1$, 故由上式推知, 当 k 时 x_k
 C , 即迭代过程恒收敛.

例 8 求 $\sqrt{115}$.

解 取初值 $x_0 = 10$, 对 $C =$

115 按(4.5)式迭代 3 次便得到精度为 10^{-6} 的结果(见表 7-6).

由于公式(4.5)对任意初值 $x_0 > 0$ 均收敛, 并且收敛的速度很快, 因此我们可取确定的初值如 $x_0 = 1$ 编制通用程序. 用这个通用程序求 $\sqrt{115}$, 也只要迭代 7 次便得到了上面的结果 10.723805.

表 7-6 计算结果

k	x_k
0	10
1	10.750000
2	10.723837
3	10.723805
4	10.723805

7.4.3 简化牛顿法与牛顿下山法

牛顿法的优点是收敛快, 缺点一是每步迭代要计算 $f(x_k)$ 及

$f(x_k)$, 计算量较大且有时 $f(x_k)$ 计算较困难, 二是初始近似 x_0 只在根 x^* 附近才能保证收敛, 如 x_0 给的不合适可能不收敛. 为克服这两个缺点, 通常可用下述方法.

(1) 简化牛顿法, 也称平行弦法. 其迭代公式为

$$x_{k+1} = x_k - C f(x_k) \quad C > 0, \quad k = 0, 1, \dots. \quad (4.7)$$

迭代函数 $\varphi(x) = x - C f(x)$.

若 $|\varphi'(x)| = |1 - C f'(x)| < 1$, 即取 $0 < C f'(x) < 2$. 在根 x^* 附近成立, 则迭代法(4.7)局部收敛.

在(4.7)中取 $C = \frac{1}{f(x_0)}$, 则

称为简化牛顿法, 这类方法计算量省, 但只有线性收敛, 其几何意义是用平行弦与 x 轴交点作为 x^* 的近似. 如图 7-4 所示.

(2) 牛顿下山法. 牛顿法收敛性依赖初值 x_0 的选取. 如果 x_0 偏离所求根 x^* 较远, 则牛顿法可能发散.

例如, 用牛顿法求解方程

$$x^3 - x - 1 = 0. \quad (4.8)$$

此方程在 $x = 1.5$ 附近的一个根

x^* . 设取迭代初值 $x_0 = 1.5$, 用牛顿法公式

$$x_{k+1} = x_k - \frac{x_k^3 - x_k - 1}{3x_k^2 - 1} \quad (4.9)$$

计算得

$$x_1 = 1.34783, \quad x_2 = 1.32520, \quad x_3 = 1.32472.$$

迭代 3 次得到的结果 x_3 有 6 位有效数字.

但是, 如果改用 $x_0 = 0.6$ 作为迭代初值, 则依牛顿法公式

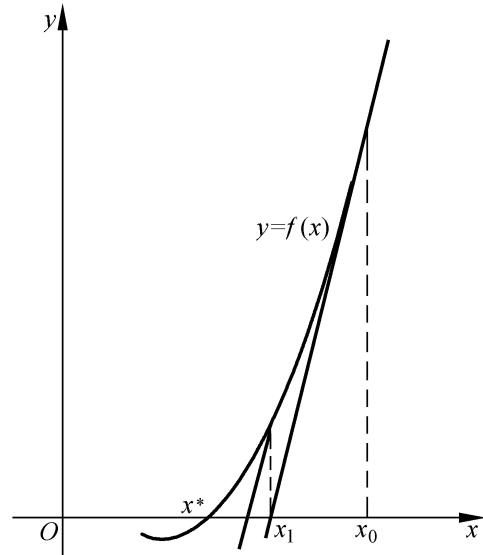


图 7-4

(4.9)迭代一次得

$$x_1 = 17.9.$$

这个结果反而比 $x_0 = 0.6$ 更偏离了所求的根 $x^* = 1.32472$.

为了防止迭代发散, 我们对迭代过程再附加一项要求, 即具有单调性:

$$|f(x_{k+1})| < |f(x_k)|. \quad (4.10)$$

满足这项要求的算法称下山法.

我们将牛顿法与下山法结合起来使用, 即在下山法保证函数值稳定下降的前提下, 用牛顿法加快收敛速度. 为此, 我们将牛顿法的计算结果

$$\bar{x}_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$$

与前一步的近似值 x_k 适当加权平均作为新的改进值

$$x_{k+1} = \bar{x}_{k+1} + (1 - \gamma)x_k, \quad (4.11)$$

其中 $(0 < \gamma < 1)$ 称为下山因子, (4.11) 即为

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} \quad (k = 0, 1, \dots), \quad (4.12)$$

称为牛顿下山法. 选择下山因子时从 $\gamma = 1$ 开始, 逐次将 γ 减半进行试算, 直到能使下降条件 (4.10) 成立为止. 若用此法解方程 (4.8), 当 $x_0 = 0.6$ 时由 (4.9) 求得 $x_1 = 17.9$, 它不满足条件 (4.10), 通过 γ 逐次取半进行试算, 当 $\gamma = \frac{1}{32}$ 时可求得

$x_1 = 1.140625$. 此时有 $f(x_1) = -0.656643$, 而 $f(x_0) = -1.384$, 显然 $|f(x_1)| < |f(x_0)|$. 由 x_1 计算 x_2, x_3, \dots 时 $\gamma = 1$, 均能使条件 (4.10) 成立. 计算结果如下:

$$x_2 = 1.36181, \quad f(x_2) = 0.1866;$$

$$x_3 = 1.32628, \quad f(x_3) = 0.00667;$$

$$x_4 = 1.32472, \quad f(x_4) = 0.0000086.$$

x_4 即为 x^* 的近似. 一般情况只要能使条件(4.10)成立, 则可得到 $\lim_k f(x_k) = 0$, 从而使 $\{x_k\}$ 收敛.

7.4.4 重根情形

设 $f(x) = (x - x^*)^m g(x)$, 整数 $m \geq 2$, $g(x^*) \neq 0$, 则 x^* 为方程 $f(x) = 0$ 的 m 重根, 此时有

$$f(x^*) = f'(x^*) = \dots = f^{(m-1)}(x^*) = 0, \quad f^{(m)}(x^*) \neq 0.$$

只要 $f'(x_k) \neq 0$ 仍可用牛顿法(4.2)计算, 此时迭代函数 $\varphi(x) = x - \frac{f(x)}{f'(x)}$ 的导数为 $\varphi'(x^*) = 1 - \frac{1}{m} \neq 0$, 且 $|\varphi'(x^*)| < 1$, 所以牛顿法求重根只是线性收敛. 若取

$$\psi(x) = x - m \frac{f(x)}{f'(x)},$$

则 $\psi(x^*) = 0$. 用迭代法

$$x_{k+1} = x_k - m \frac{f(x_k)}{f'(x_k)} \quad (k = 0, 1, \dots) \quad (4.13)$$

求 m 重根, 则具有 2 阶收敛, 但要知道 x^* 的重数 m .

构造求重根的迭代法, 还可令 $\mu(x) = f(x)/f'(x)$, 若 x^* 是 $f(x) = 0$ 的 m 重根, 则

$$\mu(x) = \frac{(x - x^*) g(x)}{mg(x) + (x - x^*) g'(x)},$$

故 x^* 是 $\mu(x) = 0$ 的单根. 对它用牛顿法, 其迭代函数为

$$\varphi(x) = x - \frac{\mu(x)}{\mu'(x)} = x - \frac{f(x)f'(x)}{[f(x)]^2 - f'(x)f(x)}.$$

从而可构造迭代法

$$x_{k+1} = x_k - \frac{f(x_k)f'(x_k)}{[f(x_k)]^2 - f'(x_k)f(x_k)} \quad (k = 0, 1, \dots), \quad (4.14)$$

它是二阶收敛的.

例 9 方程 $x^4 - 4x^2 + 4 = 0$ 的根 $x^* = 2$ 是二重根, 用上述三种方法求根.

解 先求出三种方法的迭代公式:

$$(1) \text{牛顿法} \quad x_{k+1} = x_k - \frac{x_k^2 - 2}{4x_k}.$$

$$(2) \text{用(4.13)式} \quad x_{k+1} = x_k - \frac{x_k^2 - 2}{2x_k}.$$

$$(3) \text{用(4.14)式} \quad x_{k+1} = x_k - \frac{x_k(x_k^2 - 2)}{x_k^2 + 2}.$$

取初值 $x_0 = 1.5$, 计算结果如表 7-7.

表 7-7 三种方法数值结果

k	x_k	方法(1)	方法(2)	方法(3)
1	x_1	1.458333333	1.416666667	1.411764706
2	x_2	1.436607143	1.414215686	1.414211438
3	x_3	1.425497619	1.414213562	1.414213562

计算三步, 方法(2)及(3)均达到 10 位有效数字, 而用牛顿法只有线性收敛, 要达到同样精度需迭代 30 次.

7.5 弦截法与抛物线法

用牛顿法求方程 (1.1) 的根, 每步除计算 $f(x_k)$ 外还要算 $f'(x_k)$, 当函数 $f(x)$ 比较复杂时, 计算 $f'(x)$ 往往较困难, 为此可以利用已求函数值 $f(x_k), f(x_{k-1}), \dots$ 来回避导数值 $f'(x_k)$ 的计算. 这类方法是建立在插值原理基础上的, 下面介绍两种常用方法.

7.5.1 弦截法

设 x_k, x_{k-1} 是 $f(x) = 0$ 的近似根, 我们利用 $f(x_k), f(x_{k-1})$ 构

造一次插值多项式 $p_1(x)$, 并用 $p_1(x) = 0$ 的根作为 $f(x) = 0$ 的新的近似根 x_{k+1} . 由于

$$p_1(x) = f(x_k) + \frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}}(x - x_k). \quad (5.1)$$

因此有

$$x_{k+1} = x_k - \frac{f(x_k)}{f(x_k) - f(x_{k-1})}(x_k - x_{k-1}). \quad (5.2)$$

这样导出的迭代公式(5.2)可以看做牛顿公式

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

中的导数 $f'(x_k)$ 用差商 $\frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}}$ 取代的结果.

现在解释这种迭代过程的几何意义. 如图 7-5, 曲线 $y = f(x)$ 上横坐标为 x_k, x_{k-1} 的点分别记为 P_k, P_{k-1} , 则弦线 $\overline{P_k P_{k-1}}$ 的斜率等于差商值 $\frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}}$, 其方程是

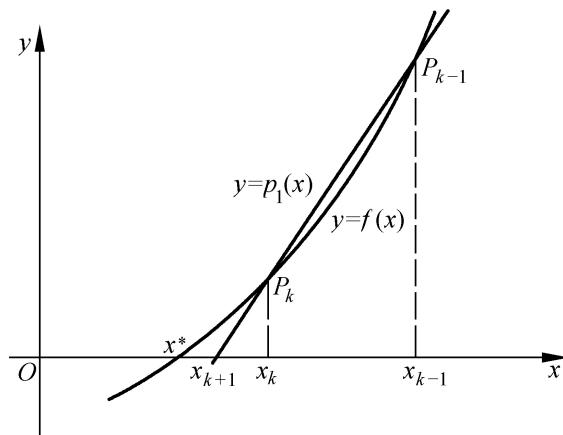


图 7-5

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

因之, 按(5.2)式求得的 x_{k+1} 实际上是弦线 $\overline{P_k P_{k-1}}$ 与 x 轴交点的横坐标. 这种算法因此而称为弦截法.

弦截法与切线法(牛顿法)都是线性化方法,但两者有本质的区别.切线法在计算 x_{k+1} 时只用到前一步的值 x_k ,而弦截法(5.2),在求 x_{k+1} 时要用到前面两步的结果 x_k, x_{k-1} ,因此使用这种方法必须先给出两个开始值 x_0, x_1 .

例 10 用弦截法解方程

$$f(x) = xe^x - 1 = 0.$$

解 设取 $x_0 = 0.5, x_1 = 0.6$ 作为开始值,用弦截法求得的结果见表 7-8,比较例 7 牛顿法的计算结果可以看出,弦截法的收敛速度也是相当快的.

实际上,下述定理断言,弦截法具有超线性的收敛性.

定理 6 假设 $f(x)$ 在根 x^* 的邻域 $:|x - x^*|$ 内具有二阶连续导数,且对任意 x 有 $f'(x) \neq 0$, 又初值 x_0, x_1 , 那么当邻域充分小时,弦截法(5.2)将按阶 $p = \frac{1 + 5}{2} = 1.618$ 收敛到根 x^* .这里 p 是方程 $x^2 - x - 1 = 0$ 的正根.

定理证明可见[2].

7.5.2 抛物线法

设已知方程 $f(x) = 0$ 的三个近似根 x_k, x_{k-1}, x_{k-2} , 我们以这三点为节点构造二次插值多项式 $p_2(x)$, 并适当选取 $p_2(x)$ 的一个零点 x_{k+1} 作为新的近似根,这样确定的迭代过程称抛物线法,亦称密勒(Müller)法.在几何图形上,这种方法的基本思想是用抛物线 $y = p_2(x)$ 与 x 轴的交点 x_{k+1} 作为所求根 x^* 的近似位置(图 7-6).

现在推导抛物线法的计算公式.插值多项式

表 7-8 计算结果

k	x_k
0	0.5
1	0.6
2	0.56532
3	0.56709
4	0.56714

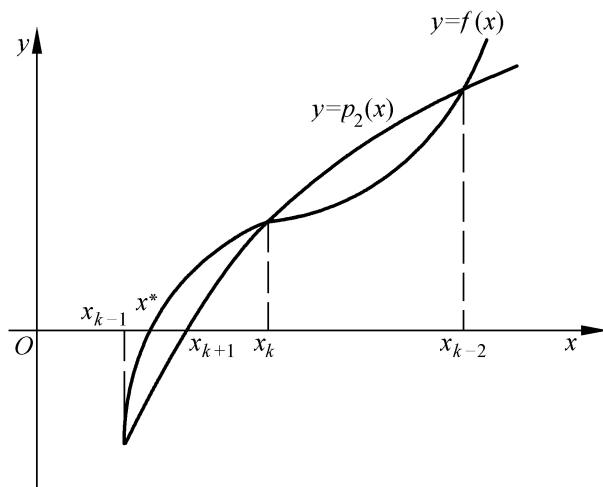


图 7-6

$$\begin{aligned} p_2(x) &= f(x_k) + f[x_k, x_{k-1}](x - x_k) \\ &\quad + f[x_k, x_{k-1}, x_{k-2}](x - x_k)(x - x_{k-1}). \end{aligned}$$

有两个零点:

$$x_{k+1} = x_k - \frac{2f(x_k)}{\pm \sqrt{4f(x_k)f[x_k, x_{k-1}, x_{k-2}]}} \quad (5.3)$$

式中

$$= f[x_k, x_{k-1}] + f[x_k, x_{k-1}, x_{k-2}](x_k - x_{k-1}).$$

为了从(5.3)式定出一个值 x_{k+1} , 我们需要讨论根式前正负号的取舍问题.

在 x_k, x_{k-1}, x_{k-2} 三个近似根中, 自然假定 x_k 更接近所求的根 x^* , 这时, 为了保证精度, 我们选式(5.3)中较接近 x_k 的一个值作为新的近似根 x_{k+1} . 为此, 只要取根式前的符号与 x^* 的符号相同.

例 11 用抛物线法求解方程 $f(x) = xe^x - 1 = 0$.

解 设用表 7-8 的前三个值

$$x_0 = 0.5, \quad x_1 = 0.6, \quad x_2 = 0.56532$$

作为开始值, 计算得

$$\begin{aligned} f(x_0) &= -0.175639, \quad f(x_1) = -0.093271, \\ f(x_2) &= -0.005031, \end{aligned}$$

$$f[x_1, x_0] = 2.68910, \quad f[x_2, x_1] = 2.83373,$$

$$f[x_2, x_1, x_0] = 2.21418.$$

故

$$= f[x_2, x_1] + f[x_2, x_1, x_0](x_2 - x_1) = 2.75694.$$

代入(5.3)式求得

$$x_3 = x_2 - \frac{2 f(x_2)}{+^2 - 4 f(x_2) f(x_2, x_1, x_0)} = 0.56714.$$

以上计算表明, 抛物线法比弦截法收敛得更快.

事实上, 在一定条件下可以证明, 对于抛物线法, 迭代误差有下列渐近关系式

$$\frac{|e_{k+1}|}{|e_k|^{1.840}} \quad \left| \frac{f(x^*)}{6 f'(x^*)} \right|^{0.42}.$$

可见抛物线法也是超线性收敛的, 其收敛的阶 $p = 1.840$ (是方程 $x^3 - x^2 - \dots - 1 = 0$ 的根), 收敛速度比弦截法更接近于牛顿法.

从(5.3)看到, 即使 x_{k-2}, x_{k-1}, x_k 均为实数, x_{k+1} 也可以是复数, 所以抛物线法适用于求多项式的实根和复根.

7.6 解非线性方程组的牛顿迭代法

考虑方程组

$$\begin{aligned} f_1(x_1, \dots, x_n) &= 0, \\ &\dots \\ f_n(x_1, \dots, x_n) &= 0. \end{aligned} \tag{6.1}$$

其中 f_1, \dots, f_n 均为 (x_1, \dots, x_n) 的多元函数. 若用向量记号记 $\mathbf{x} = (x_1, \dots, x_n)^T \in \mathbf{R}^n$, $\mathbf{F} = (f_1, \dots, f_n)^T$, (6.1)就可写成

$$\mathbf{F}(\mathbf{x}) = \mathbf{0}. \tag{6.2}$$

当 $n \geq 2$, 且 $f_i (i=1, \dots, n)$ 中至少有一个是自变量 $x_i (i=1, \dots, n)$ 的非线性函数时, 则称方程组(6.1)为非线性方程组. 非线性方程

组求根问题是前面介绍的方程(即 $n=1$)求根的直接推广, 实际上只要把前面介绍的单变量函数 $f(x)$ 看成向量函数 $\mathbf{F}(\mathbf{x})$ 则可将单变量方程求根方法推广到方程组(6.2). 若已给出方程(6.2)的一个近似根 $\mathbf{x}^{(k)} = (x_1^{(k)}, \dots, x_n^{(k)})^T$, 将函数 $\mathbf{F}(\mathbf{x})$ 的分量 $f_i(\mathbf{x})$ ($i=1, \dots, n$) 在 $\mathbf{x}^{(k)}$ 用多元函数泰勒展开, 并取其线性部分, 则可表示为

$$\mathbf{F}(\mathbf{x}) - \mathbf{F}(\mathbf{x}^{(k)}) + \mathbf{F}'(\mathbf{x}^{(k)})(\mathbf{x} - \mathbf{x}^{(k)}).$$

令上式右端为零, 得到线性方程组

$$\mathbf{F}'(\mathbf{x}^{(k)})(\mathbf{x} - \mathbf{x}^{(k)}) = -\mathbf{F}(\mathbf{x}^{(k)}), \quad (6.3)$$

其中

$$\begin{aligned} \mathbf{F}'(\mathbf{x}) = & \begin{array}{cccc} \frac{f_1(\mathbf{x})}{x_1} & \frac{f_1(\mathbf{x})}{x_2} & \cdots & \frac{f_1(\mathbf{x})}{x_n} \\ & x_1 & x_2 & \cdots & x_n \\ \frac{f_2(\mathbf{x})}{x_1} & \frac{f_2(\mathbf{x})}{x_2} & \cdots & \frac{f_2(\mathbf{x})}{x_n} \\ x_1 & x_2 & \cdots & x_n \\ \cdots & \cdots & \cdots & \cdots \\ \frac{f_n(\mathbf{x})}{x_1} & \frac{f_n(\mathbf{x})}{x_2} & \cdots & \frac{f_n(\mathbf{x})}{x_n} \\ x_1 & x_2 & \cdots & x_n \end{array} \end{aligned} \quad (6.4)$$

称为 $\mathbf{F}(\mathbf{x})$ 的雅可比 (Jacobi) 矩阵. 求解线性方程组(6.3), 并记解为 $\mathbf{x}^{(k+1)}$, 则得

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{F}'(\mathbf{x}^{(k)})^{-1} \mathbf{F}(\mathbf{x}^{(k)}) \quad (k=0, 1, \dots). \quad (6.5)$$

这就是解非线性方程组(6.2)的牛顿迭代法.

例 12 求解方程组

$$f_1(x_1, x_2) = x_1 + 2x_2 - 3 = 0,$$

$$f_2(x_1, x_2) = 2x_1^2 + x_2^2 - 5 = 0.$$

给定初值 $\mathbf{x}^{(0)} = (1.5, 1.0)^T$, 用牛顿法求解.

解 先求雅可比矩阵

$$\mathbf{F}'(\mathbf{x}) = \begin{pmatrix} 1 & 2 \\ 4x_1 & 2x_2 \end{pmatrix}, \quad \mathbf{F}'(\mathbf{x})^{-1} = \begin{pmatrix} 1 & 2x_2 \\ 2x_2 & 8x_1 \end{pmatrix}^{-1} = \begin{pmatrix} 1 & 2x_2 \\ 2x_2 & 8x_1 \end{pmatrix}^{-1} = \begin{pmatrix} 1 & 2 \\ 2 & 4x_1 \end{pmatrix}.$$

由牛顿法(6.5)得

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \frac{1}{2x_2^{(k)} - 8x_1^{(k)}} \begin{pmatrix} 2x_2^{(k)} - 2 & x_1^{(k)} + 2x_2^{(k)} - 3 \\ -4x_1^{(k)} - 1 & 2(x_1^{(k)})^2 + (x_2^{(k)})^2 - 5 \end{pmatrix},$$

即

$$x_1^{(k+1)} = x_1^{(k)} - \frac{(x_2^{(k)})^2 - 2(x_1^{(k)})^2 + x_1^{(k)}x_2^{(k)} - 3x_2^{(k)} + 5}{x_2^{(k)} - 4x_1^{(k)}},$$

$$x_2^{(k+1)} = x_2^{(k)} - \frac{(x_2^{(k)})^2 - 2(x_1^{(k)})^2 - 8x_1^{(k)}x_2^{(k)} + 12x_2^{(k)} - 5}{2(x_2^{(k)} - 4x_1^{(k)})},$$

$$(k = 0, 1, \dots).$$

由 $\mathbf{x}^{(0)} = (1.5, 1.0)^T$ 逐次迭代得到

$$\mathbf{x}^{(1)} = (1.5, 0.75)^T,$$

$$\mathbf{x}^{(2)} = (1.488095, 0.755952)^T,$$

$$\mathbf{x}^{(3)} = (1.488034, 0.755983)^T.$$

$\mathbf{x}^{(3)}$ 的每一位都是有效数字 .

评注

本章着重介绍求解单变量非线性方程 $f(x) = 0$ 的迭代法及其理论. 不动点迭代、局部收敛性及收敛阶等基本概念是十分重要的, 它很容易推广到非线性方程组. 在迭代法中以牛顿法最实用, 它在单根附近具有 2 阶收敛, 但应用时要选取较好的初始近似才能保证迭代收敛. 为克服这一缺点, 可使用牛顿下山法. 斯蒂芬森方法可将一阶方法加速为二阶, 也是值得重视的算法. 弦截法(或称割线法)与抛物线法(也称密勒法)是属于插值方法, 它们不用算 $f(x)$ 的导数, 又具有超线性收敛, 也是常用的有效方法. 这类方法是多点迭代法, 它不同于 $x_{k+1} = (x_k)$ 的单点迭代, 计算时必须给出两个以上的初始近似. 其收敛性说明可参看文献[2]和[7].

非线性方程组的解法和理论是当今数值分析研究的重要课题之一, 新方法不断出现, 它也是科学计算经常遇到的. 这里我们只

简单介绍了牛顿法的迭代公式,有关理论均未提及,需要进一步了解的读者可参阅文献[1]和[17].

单个代数方程(即多项式方程)求根有久远的历史,也有不少特殊方法.本章均未介绍,有关理论和算法读者可参阅文献[7]和[18]及其他文献.

习 题

1. 用二分法求方程 $x^2 - x - 1 = 0$ 的正根,要求误差小于 0.05.

2. 为求方程 $x^3 - x^2 - 1 = 0$ 在 $x_0 = 1.5$ 附近的一个根,设将方程改写成下列等价形式,并建立相应的迭代公式.

(1) $x = 1 + 1/x^2$, 迭代公式 $x_{k+1} = 1 + 1/x_k^2$;

(2) $x^3 = 1 + x^2$, 迭代公式 $x_{k+1} = \sqrt[3]{1 + x_k^2}$;

(3) $x^2 = \frac{1}{x - 1}$, 迭代公式 $x_{k+1} = \sqrt{x_k - 1}$.

试分析每种迭代公式的收敛性,并选取一种公式求出具有四位有效数字的近似根.

3. 比较求 $e^x + 10x - 2 = 0$ 的根到三位小数所需的计算量:

(1) 在区间 $[0, 1]$ 内用二分法;

(2) 用迭代法 $x_{k+1} = (2 - e^{x_k})/10$, 取初值 $x_0 = 0$.

4. 给定函数 $f(x)$, 设对一切 x , $f'(x)$ 存在且 $0 < m < f'(x) < M$, 证明对于范围 $0 < |x| < 2/M$ 内的任意定数 ϵ , 迭代过程 $x_{k+1} = x_k - f(x_k)/f'(x_k)$ 均收敛于 $f(x) = 0$ 的根 x^* .

5. 用斯蒂芬森迭代法计算第 2 题中(2),(3)的近似根,精确到 10^{-5} .

6. 设 $\varphi(x) = x - p(x)f(x) - q(x)f^2(x)$, 试确定函数 $p(x)$ 和 $q(x)$, 使求解 $f(x) = 0$ 且以 $\varphi(x)$ 为迭代函数的迭代法至少三阶收敛.

7. 用下列方法求 $f(x) = x^3 - 3x - 1 = 0$ 在 $x_0 = 2$ 附近的根. 根的准确值 $x^* = 1.87938524\dots$, 要求计算结果准确到四位有效数字.

(1) 用牛顿法;

(2) 用弦截法,取 $x_0 = 2$, $x_1 = 1.9$;

(3) 用抛物线法, 取 $x_0 = 1, x_1 = 3, x_2 = 2$.

8. 分别用二分法和牛顿法求 $x - \tan x = 0$ 的最小正根.

9. 研究求 a 的牛顿公式

$$x_{k+1} = \frac{1}{2} x_k + \frac{a}{x_k}, \quad x_0 > 0.$$

证明对一切 $k = 1, 2, \dots, x_k = a$ 且序列 x_1, x_2, \dots 是递减的.

10. 对于 $f(x) = 0$ 的牛顿公式 $x_{k+1} = x_k - f(x_k)/f'(x_k)$, 证明

$$R_k = (x_k - x_{k-1}) / (x_{k-1} - x_{k-2})^2$$

收敛到 $-f'(x^*)/[2f''(x^*)]$, 这里 x^* 为 $f(x) = 0$ 的根.

11. 用牛顿法和求重根迭代法 (4.13) 和 (4.14) 计算方程 $f(x) = \sin x - \frac{x^2}{2} = 0$ 的一个近似根, 准确到 10^{-5} , 初始值 $x_0 = \frac{\pi}{2}$.

12. 应用牛顿法于方程 $x^3 - a = 0$, 导出求立方根 a 的迭代公式, 并讨论其收敛性.

13. 应用牛顿法于方程 $f(x) = 1 - \frac{a}{x^2} = 0$, 导出求 a 的迭代公式, 并用此

公式求 $\sqrt[11]{5}$ 的值.

14. 应用牛顿法于方程 $f(x) = x^n - a = 0$ 和 $f(x) = 1 - \frac{a}{x^n} = 0$, 分别导出求 $\sqrt[n]{a}$ 的迭代公式, 并求

$$\lim_k \frac{x_k^n - a}{x_{k+1}^n} / \frac{x_k^n - a}{x_k} = 0.$$

15. 证明迭代公式

$$x_{k+1} = \frac{x_k(x_k^2 + 3a)}{3x_k^2 + a}$$

是计算 a 的三阶方法. 假定初值 x_0 充分靠近根 x^* , 求

$$\lim_x \frac{a - x_{k+1}}{a - x_k} / \frac{a - x_k}{x_k} = 0.$$

16. 用牛顿法解方程组

$$x^2 + y^2 = 4,$$

$$x^2 - y^2 = 1.$$

取 $\mathbf{x}^0 = (1, 1, 2)^T$.

第8章 矩阵特征值问题计算

8.1 引言

物理、力学和工程技术中的很多问题在数学上都归结为求矩阵的特征值问题。例如，振动问题（大型桥梁或建筑物的振动、机械的振动、电磁振荡等），物理学中某些临界值的确定，这些问题都归结为下述数学问题。

定义 1 (1) 已知 $\mathbf{A} = (a_{ij})_{n \times n}$, 则称

$$\begin{aligned} & - a_{11} \quad - a_{12} \quad \dots \quad - a_{1n} \\ (\lambda) = \det(\mathbf{I} - \mathbf{A}) = \det & \begin{array}{cccc} - a_{11} & - a_{12} & \dots & - a_{1n} \\ \cdots & \cdots & \text{W} & \cdots \\ - a_{n1} & - a_{n2} & \dots & - a_{nn} \end{array} \\ & = (-1)^n (a_{11} + a_{22} + \dots + a_{nn})^{n-1} + (\text{次级 } n-2 \text{ 的项}) \end{aligned}$$

为 \mathbf{A} 的特征多项式。

A 的特征方程

$$(\lambda) = \det(\mathbf{I} - \mathbf{A}) = 0 \quad (1.1)$$

一般有 n 个根（实的或复的，重根按重数计算）（当 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 时， $(\lambda) = 0$ 为实系数 n 次代数方程，其复根共轭成对出现），称为 \mathbf{A} 的特征值。用 (\mathbf{A}) 表示 \mathbf{A} 的所有特征值的集合。

(2) 设 λ 为 \mathbf{A} 特征值，相应的齐次方程组

$$(\mathbf{I} - \mathbf{A})\mathbf{x} = \mathbf{0} \quad (1.2)$$

的非零解 \mathbf{x} 称为矩阵 \mathbf{A} 的对应于 λ 的特征向量。

例 1 求 \mathbf{A} 的特征值及特征向量，其中

$$\mathbf{A} = \begin{pmatrix} 2 & 1 & 0 \\ 1 & 3 & 1 \\ 0 & 1 & 2 \end{pmatrix}.$$

解 矩阵 \mathbf{A} 的特征方程为

$$(\lambda) = \det(\mathbf{I} - \mathbf{A}) = \lambda^3 - 7\lambda^2 + 14\lambda - 8 = 0,$$

求得 \mathbf{A} 特征值为：

$$\lambda_1 = 1, \lambda_2 = 2, \lambda_3 = 4.$$

对应于各特征值矩阵 \mathbf{A} 特征向量分别为：

$$\mathbf{x}_1 = \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix}, \quad \mathbf{x}_2 = \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix}, \quad \mathbf{x}_3 = \begin{pmatrix} 1 \\ 2 \\ 1 \end{pmatrix}.$$

下面叙述有关特征值的一些结果。

定理 1 设 λ 为 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 的特征值且 $\mathbf{Ax} = \lambda \mathbf{x}$, 其中 $\mathbf{x} \neq \mathbf{0}$, 则

(1) $c\lambda$ 为 $c\mathbf{A}$ 的特征值 (c 为常数 $c \neq 0$);

(2) $\lambda - p$ 为 $\mathbf{A} - p\mathbf{I}$ 的特征值, 即 $(\mathbf{A} - p\mathbf{I})\mathbf{x} = (\lambda - p)\mathbf{x}$;

(3) λ^k 为 \mathbf{A}^k 的特征值;

(4) 设 \mathbf{A} 为非奇异阵, 那么 0 且 $\frac{1}{\lambda}$ 为 \mathbf{A}^{-1} 特征值, 即

$$\mathbf{A}^{-1}\mathbf{x} = \frac{1}{\lambda}\mathbf{x}.$$

定理 2 设 a_{ii} ($i = 1, 2, \dots, n$) 为 n 阶矩阵 $\mathbf{A} = (a_{ij})_{n \times n}$ 特征值, 则

$$(1) \sum_{i=1}^n a_{ii} = \text{tr}(\mathbf{A});$$

$$(2) \det(\mathbf{A}) = a_{11}a_{22} \dots a_{nn}.$$

定理 3 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$, 则

$$(\mathbf{A}^T) = (\mathbf{A}).$$

定理 4 设 \mathbf{A} 为分块上三角阵, 即

$$\mathbf{A} = \begin{matrix} \mathbf{A}_1 & \mathbf{A}_2 & \dots & \mathbf{A}_m \\ & \mathbf{A}_2 & \dots & \mathbf{A}_m \\ & & \ddots & \dots \\ & & & \mathbf{A}_{m-m} \end{matrix},$$

其中每个对角块 \mathbf{A}_i 均为方阵, 则 $(\mathbf{A}) = \begin{matrix} m \\ i=1 \end{matrix} (\mathbf{A}_i)$.

定理 5 设 \mathbf{A} 与 \mathbf{B} 为相似矩阵 (即存在非奇异阵 \mathbf{P} 使 $\mathbf{B} = \mathbf{P}^{-1} \mathbf{AP}$), 则

(1) \mathbf{A} 与 \mathbf{B} 有相同的特征值;

(2) 如果 \mathbf{y} 是 \mathbf{B} 特征向量, 则 \mathbf{Py} 是 \mathbf{A} 特征向量.

定理 5 说明, 一个矩阵 \mathbf{A} 经过相似变换 ($\mathbf{A} \rightarrow \mathbf{B} = \mathbf{P}^{-1} \mathbf{AP}$), 则 \mathbf{A} 的特征值不变.

定义 2 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$, 如果 \mathbf{A} 有一个重数为 k 的特征值 且 对应于 矩阵 \mathbf{A} 的线性无关的特征向量个数少于 k (一般 $< k$), 称 \mathbf{A} 为亏损矩阵.

一个亏损矩阵是一个没有足够特征向量的矩阵, 亏损矩阵在理论上和计算上都存在困难.

定理 6 (1) $\mathbf{A} \in \mathbb{R}^{n \times n}$ 可对角化, 即存在非奇异矩阵 \mathbf{P} 使

$$\mathbf{P}^{-1} \mathbf{AP} = \begin{matrix} 1 & & & \\ & 2 & & \\ & & \ddots & \\ & & & n \end{matrix}$$

的充要条件是 \mathbf{A} 具有 n 个线性无关的特征向量.

(2) 如果 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 有 m 个 ($m < n$) 不同的特征值 $\lambda_1, \lambda_2, \dots, \lambda_m$, 则对应的特征向量 $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m$ 线性无关.

定理 7 (对称矩阵的正交约化) 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 为对称矩阵, 则:

(1) \mathbf{A} 的特征值均为实数;

(2) \mathbf{A} 有 n 个线性无关的特征向量;

(3) 存在一个正交矩阵 \mathbf{P} 使得

$$\mathbf{P}^T \mathbf{A} \mathbf{P} = \begin{matrix} & 1 \\ & & 2 \\ & & & \ddots \\ & & & & n \\ & & & & & \mathbf{W} \end{matrix},$$

且 λ_i ($i = 1, \dots, n$) 为 \mathbf{A} 特征值, 而 $\mathbf{P} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n)$ 的列向量 \mathbf{u}_i 为 \mathbf{A} 的对应于 λ_i 的特征向量.

下面讨论矩阵特征值界的估计.

定义 3 设 $\mathbf{A} = (a_{ij})_{n \times n}$. 令: (1) $r_i = \max_{\substack{j=1 \\ j \neq i}} |a_{ij}|$ ($i = 1, 2, \dots, n$); (2) 集合 $D_i = \{z / |z - a_{ii}| \leq r_i, z \in \mathbf{C}\}$. 称复平面上以 a_{ii} 为圆心, 以 r_i 为半径的所有圆盘为 \mathbf{A} 的 Gershgorin 圆盘.

定理 8(Gershgorin 圆盘定理) (1) 设 $\mathbf{A} = (a_{ij})_{n \times n}$, 则 \mathbf{A} 的每一个特征值必属于下述某个圆盘之中

$$|z - a_{ii}| \leq r_i = \max_{\substack{j=1 \\ j \neq i}} |a_{ij}| \quad (i = 1, 2, \dots, n).$$

或者说, \mathbf{A} 的特征值都在复平面上 n 个圆盘的并集中.

(2) 如果 \mathbf{A} 有 m 个圆盘组成一个连通的并集 S , 且 S 与余下 $n - m$ 个圆盘是分离的, 则 S 内恰包含 \mathbf{A} 的 m 个特征值.

特别地, 如果 \mathbf{A} 的一个圆盘 D_i 是与其他圆盘分离的(即孤立圆盘), 则 D_i 中精确地包含 \mathbf{A} 的一个特征值.

证明 只就(1)给出证明. 设 λ 为 \mathbf{A} 的特征值, 即

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{x}, \text{ 其中 } \mathbf{x} = (x_1, x_2, \dots, x_n)^T \neq \mathbf{0}.$$

记 $|x_k| = \max_{i=1}^n |x_i| = \| \mathbf{x} \|$. 考虑 $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ 的第 k 个方程, 即

n

$$\sum_{j=1}^n a_{kj} x_j = x_k,$$

或

$$(- a_{kk}) x_k = \sum_{j \neq k} a_{kj} x_j,$$

于是 $| - a_{kk} | / | x_k | = \sum_{j \neq k} | a_{kj} | / | x_j | = | x_k | / | x_k | = 1$,

即

$$| - a_{kk} | / \sum_{j \neq k} | a_{kj} | = n.$$

这说明, \mathbf{A} 的每一个特征值必位于 \mathbf{A} 的一个圆盘中, 并且相应的特征值一定位于第 k 个圆盘中(其中 k 是对应特征向量 \mathbf{x} 绝对值最大的分量的下标).

利用相似矩阵性质, 有时可以获得 \mathbf{A} 的特征值进一步的估计, 即适当选取非奇异对角阵

$$\mathbf{D}^{-1} = \begin{matrix} & & 1 \\ & & 2 \\ & & \ddots \\ 1 & & & n \end{matrix}$$

并做相似变换 $\mathbf{D}^{-1} \mathbf{A} \mathbf{D} = \begin{matrix} a_{ij} \\ i \\ n \times n \end{matrix}$. 适当选取 ω_i ($i = 1, 2, \dots, n$) 可

使某些圆盘半径及连通性发生变化.

例 2 估计矩阵

$$\mathbf{A} = \begin{matrix} 4 & 1 & 0 \\ 1 & 0 & -1 \\ 1 & 1 & -4 \end{matrix}$$

特征值的范围.

解 \mathbf{A} 的 3 个圆盘为

$$D_1 : | -4 | = 1,$$

$$D_2 : | / | = 2,$$

$$D_3 : | +4 | = 2.$$

由定理 8, 可知 \mathbf{A} 的 3 个特征值位于 3 个圆盘的并集中, 由于 D_1 是孤立圆盘, 所以 D_1 内恰好包含 \mathbf{A} 的一个特征值 λ_1 (为实特征值), 即

$$3 \quad -1 \quad 5$$

\mathbf{A} 的其他两个特征值 λ_2, λ_3 包含在 D_2, D_3 的并集中.

现选取对角阵

$$\mathbf{D}^{-1} = \begin{matrix} 1 \\ & 1 \\ & & 0.9 \end{matrix}$$

做相似变换

$$\mathbf{A}^{-1} \mathbf{A} = \mathbf{D}^{-1} \mathbf{A} \mathbf{D} = \begin{matrix} 4 & 1 & 0 \\ 1 & 0 & -\frac{10}{9} \\ 0.9 & 0.9 & -4 \end{matrix}.$$

\mathbf{A} 的 3 个圆盘为

$$E_1 : / -4 / 1,$$

$$E_2 : / / \frac{19}{9},$$

$$E_3 : / +4 / 1.8.$$

显然, 3 个圆盘都是孤立圆盘, 所以, 每一个圆盘都包含 \mathbf{A} 的一个特征值(为实特征值)且有估计

$$\begin{aligned} & 3 \quad -1 \quad 5, \\ & -\frac{19}{9} \quad 2 \quad \frac{19}{9}, \\ & -5.8 \quad 3 \quad -2.2. \end{aligned}$$

下面给出理论上有关通过酉相似变换及正交相似变换可以约化一般矩阵 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 到什么程度的问题.

定理 9(Schur 定理) 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$, 则存在酉阵 \mathbf{U} 使

$$\mathbf{U}^H \mathbf{A} \mathbf{U} = \begin{matrix} r_{11} & r_{12} & \cdots & r_{1n} \\ & r_{22} & \cdots & r_{2n} \\ & & \ddots & \cdots \\ & & & r_{nn} \end{matrix} = \mathbf{R} \text{ (上三角阵),}$$

其中 r_{ii} ($i = 1, 2, \dots, n$) 为 \mathbf{A} 的特征值 .

当 $\mathbf{A} \in \mathbf{R}^{n \times n}$ 时, 如果限制用正交相似变换, 由于 \mathbf{A} 有复的特征值, \mathbf{A} 不能用正交相似变换约化为上三角阵 . 用正交相似变换能将 $\mathbf{A} \in \mathbf{R}^{n \times n}$ 约化到什么程度呢 ?

定理 10(实 Schur 分解) 设 $\mathbf{A} \in \mathbf{R}^{n \times n}$, 则存在正交矩阵 \mathbf{Q} 使

$$\mathbf{Q}^T \mathbf{A} \mathbf{Q} = \begin{matrix} \mathbf{R}_1 & & \cdots & & \mathbf{R}_m \\ & \mathbf{R}_2 & & \cdots & & \mathbf{R}_m \\ & & \ddots & & & \cdots \\ & & & \mathbf{R}_{m-m} & & \end{matrix},$$

其中对角块 \mathbf{R}_i ($i = 1, 2, \dots, m$) 为一阶或二阶方阵, 且每个一阶 \mathbf{R}_i 是 \mathbf{A} 的实特征值, 每个二阶对角块 \mathbf{R}_j 的两个特征值是 \mathbf{A} 的两个共轭复特征值 .

我们转向实 Schur 型的实际计算 .

定义 4 设 \mathbf{A} 为 n 阶实对称矩阵, 对于任一非零向量 \mathbf{x} , 称

$$R(\mathbf{x}) = \frac{\langle \mathbf{Ax}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle}$$

为对应于向量 \mathbf{x} 的瑞利 (Rayleigh) 商 .

定理 11 设 $\mathbf{A} \in \mathbf{R}^{n \times n}$ 为对称矩阵 (其特征值次序记为 $\lambda_1, \lambda_2, \dots, \lambda_n$), 则

$$1. \quad \lambda_1 = \frac{\langle \mathbf{Ax}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} \quad (\text{对任何非零 } \mathbf{x} \in \mathbf{R}^n);$$

$$2. \quad \lambda_n = \max_{\substack{\mathbf{x} \in \mathbf{R}^n \\ \mathbf{x} \neq \mathbf{0}}} \frac{\langle \mathbf{Ax}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle};$$

$$3. \quad n = \min_{\substack{\mathbf{x} \in \mathbb{R}^n \\ \mathbf{x} \neq \mathbf{0}}} \frac{(\mathbf{Ax}, \mathbf{x})}{(\mathbf{x}, \mathbf{x})}.$$

证明 只证 1, 关于 2,3 留作习题 .

由于 \mathbf{A} 为实对称矩阵, 可将 $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ 对应的特征向量 $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ 正交规范化, 则有 $(\mathbf{x}_i, \mathbf{x}_j) = \delta_{ij}$. 设 $\mathbf{x} \neq \mathbf{0}$ 为 \mathbb{R}^n 中任一向量, 则有展开式

$$\mathbf{x} = \sum_{i=1}^n x_i \mathbf{x}_i, \quad \mathbf{x}_2 = \sum_{i=1}^n \begin{pmatrix} 2 & & \\ & i & \\ & & 2 \end{pmatrix} x_i \mathbf{x}_i = 0,$$

于是

$$\frac{(\mathbf{Ax}, \mathbf{x})}{(\mathbf{x}, \mathbf{x})} = \frac{\sum_{i=1}^n a_{ii} x_i^2}{\sum_{i=1}^n x_i^2}.$$

从而 1 成立. 结论 1 说明瑞利商必位于 λ_n 和 λ_1 之间 .

关于计算矩阵 \mathbf{A} 的特征值问题, 当 $n=2, 3$ 时, 我们还可按行列式展开的办法求 $\det(\mathbf{A} - \lambda \mathbf{I}) = 0$ 的根. 但当 n 较大时, 如果按展开行列式的办法, 首先求出 $\det(\mathbf{A} - \lambda \mathbf{I})$ 的系数, 再求 $\det(\mathbf{A} - \lambda \mathbf{I}) = 0$ 的根, 工作量就非常大, 用这种办法求矩阵特征值是不切实际的, 由此需要研究求 \mathbf{A} 的特征值及特征向量的数值解法 .

本章将介绍一些计算机上常用的两类方法, 一类是幂法及反幂法(迭代法), 另一类是正交相似变换的方法(变换法) .

8.2 幂法及反幂法

8.2.1 幂法

幂法是一种计算矩阵主特征值(矩阵按模最大的特征值)及对应特征向量的迭代方法, 特别适用于大型稀疏矩阵. 反逆法是计算海森伯格阵或三对角阵的对应一个给定近似特征值的特征向量的

有效方法之一 .

设实矩阵 $\mathbf{A} = (a_{ij})_{n \times n}$ 有一个完全的特征向量组, 其特征值为 $\lambda_1, \lambda_2, \dots, \lambda_n$, 相应的特征向量为 $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$. 已知 \mathbf{A} 的主特征值是实根, 且满足条件

$$|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq \dots \geq |\lambda_n|, \quad (2.1)$$

现讨论求 λ_1 及 \mathbf{x}_1 的方法 .

幂法的基本思想是任取一个非零的初始向量 \mathbf{v}_0 , 由矩阵 \mathbf{A} 构造一向量序列

$$\begin{aligned} \mathbf{v}_1 &= \mathbf{Av}_0 \\ \mathbf{v}_2 &= \mathbf{Av}_1 = \mathbf{A}^2 \mathbf{v}_0, \\ &\dots \\ \mathbf{v}_{k+1} &= \mathbf{Av}_k = \mathbf{A}^{k+1} \mathbf{v}_0, \\ &\dots \end{aligned} \quad (2.2)$$

称为迭代向量. 由假设, \mathbf{v}_0 可表示为

$$\mathbf{v}_0 = \lambda_1 \mathbf{x}_1 + \lambda_2 \mathbf{x}_2 + \dots + \lambda_n \mathbf{x}_n \quad (\text{设 } \lambda_1 \neq 0), \quad (2.3)$$

于是

$$\begin{aligned} \mathbf{v}_k &= \mathbf{Av}_{k-1} = \mathbf{A}^k \mathbf{v}_0 = \lambda_1^k \mathbf{x}_1 + \lambda_2^k \mathbf{x}_2 + \dots + \lambda_n^k \mathbf{x}_n \\ &= \lambda_1^k \mathbf{x}_1 + \sum_{i=2}^n i(\sqrt{\lambda_1})^k \mathbf{x}_i = \lambda_1^k (\mathbf{x}_1 + \mathbf{x}_k), \end{aligned}$$

其中 $\mathbf{x}_k = \sum_{i=2}^n i(\sqrt{\lambda_1})^k \mathbf{x}_i$. 由假设

$|\sqrt{\lambda_1}| < 1 (i = 2, 3, \dots, n)$, 故 $\lim_{k \rightarrow \infty} \mathbf{x}_k = \mathbf{0}$, 从而

$$\lim_{k \rightarrow \infty} \frac{\mathbf{v}_k}{\mathbf{x}_k} = \lambda_1 \mathbf{x}_1. \quad (2.4)$$

这说明序列 $\frac{\mathbf{v}_k}{\mathbf{x}_k}$ 越来越接近 \mathbf{A} 的对应于 λ_1 的特征向量, 或者说

当 k 充分大时

$$\mathbf{v}_k \approx \lambda_1 \mathbf{x}_1, \quad (2.5)$$

即迭代向量 \mathbf{v}_k 为 λ_1 的特征向量的近似向量(除一个因子外) .

下面再考虑主特征值 λ_1 的计算, 用 $(\mathbf{v}_k)_i$ 表示 \mathbf{v}_k 的第 i 个分量, 则

$$\frac{(\mathbf{v}_{k+1})_i}{(\mathbf{v}_k)_i} = \lambda_1 - \frac{\lambda_1 (\mathbf{x}_1)_i + (\mathbf{x}_{k+1})_i}{\lambda_1 (\mathbf{x}_1)_i + (\mathbf{x}_k)_i}, \quad (2.6)$$

故

$$\lim_{k \rightarrow \infty} \frac{(\mathbf{v}_{k+1})_i}{(\mathbf{v}_k)_i} = \lambda_1, \quad (2.7)$$

也就是说两相邻迭代向量分量的比值收敛到主特征值 .

这种由已知非零向量 \mathbf{v} 及矩阵 \mathbf{A} 的乘幂 \mathbf{A}^k 构造向量序列 $\{\mathbf{v}_k\}$ 以计算 \mathbf{A} 的主特征值 λ_1 (利用(2.7)式)及相应特征向量(利用(2.5)式)的方法称为幂法 .

由(2.6)式知, $\frac{(\mathbf{v}_{k+1})_i}{(\mathbf{v}_k)_i}$ 的收敛速度由比值 $r = \left| \frac{\lambda_2}{\lambda_1} \right|$ 来确定, r 越小收敛越快, 但当 $r = \left| \frac{\lambda_2}{\lambda_1} \right| = 1$ 时收敛可能就很慢 .

总结上述讨论, 有

定理 12 设 $\mathbf{A} \in \mathbf{R}^{n \times n}$ 有 n 个线性无关的特征向量, 主特征值 λ_1 满足

$$|\lambda_1| > |\lambda_2| > |\lambda_3| > \dots > |\lambda_n|,$$

则对任何非零初始向量 $\mathbf{v}(=0)$, (2.4), (2.7)式成立 .

如果 \mathbf{A} 的主特征值为实的重根, 即 $\lambda_1 = \lambda_2 = \dots = \lambda_r$, 且

$$|\lambda_r| > |\lambda_{r+1}| > \dots > |\lambda_n|,$$

又设 \mathbf{A} 有 n 个线性无关的特征向量, λ_1 对应的 r 个线性无关特征向量为 $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_r$, 则由(2.2)式

$$\begin{aligned} \mathbf{v}_k &= \mathbf{A}^k \mathbf{v}_0 = \sum_{i=1}^k \lambda_i \mathbf{x}_i + \sum_{i=r+1}^n \lambda_i (\sqrt[r]{\lambda_1})^k \mathbf{x}_i, \\ \lim_k \frac{\mathbf{v}_k}{\lambda_1^k} &= \sum_{i=1}^r \lambda_i \mathbf{x}_i \quad \text{设 } \sum_{i=1}^r \lambda_i \mathbf{x}_i = \mathbf{0}. \end{aligned}$$

这说明当 A 的主特征值是实的重根时, 定理 5 的结论还是正确的.

应用幂法计算 A 的主特征值 λ_1 及对应的特征向量时, 如果 $|\lambda_1| > 1$ (或 $|\lambda_1| < 1$), 迭代向量 v 的各个不等于零的分量将随 k 而趋于无穷(或趋于零), 这样在计算机实现时就可能“溢出”. 为了克服这个缺点, 就需要将迭代向量加以规范化.

设有一向量 $v \neq 0$, 将其规范化得到向量

$$u = \frac{v}{\max(v)},$$

其中 $\max(v)$ 表示向量 v 的绝对值最大的分量, 即如果有

$$|v_{i_0}| = \max_i |v_i|,$$

则 $\max(v) = v_{i_0}$, 且 i_0 为所有绝对值最大的分量中的最小下标.

在定理 12 的条件下幂法可这样进行: 任取一初始向量 $v \neq 0$ ($v_1 \neq 0$), 构造向量序列 $\max(v)$

$$\begin{aligned} v &= Av = Av_0, & u &= \frac{v}{\max(v)} = \frac{Av_0}{\max(Av_0)}, \\ v &= A^2 v = \frac{A^2 v}{\max(Av_0)}, & u &= \frac{v}{\max(v)} = \frac{A^2 v}{\max(A^2 v)}, \\ &\dots &&\dots \\ v_k &= \frac{A^k v_0}{\max(A^{k-1} v_0)}, & u_k &= \frac{A^k v_0}{\max(A^k v_0)}, \end{aligned}$$

由(2.3)式

$$\begin{aligned} A^k v_0 &= \sum_{i=1}^n \lambda_i^k \mathbf{x}_i = \lambda_1^k \mathbf{x}_1 + \sum_{i=2}^n \lambda_i^k \mathbf{x}_i, \quad (2.8) \\ u_k &= \frac{A^k v_0}{\max(A^k v_0)} = \frac{\lambda_1^k \mathbf{x}_1 + \sum_{i=2}^n \lambda_i^k \mathbf{x}_i}{\max(\lambda_1^k \mathbf{x}_1 + \sum_{i=2}^n \lambda_i^k \mathbf{x}_i)} \end{aligned}$$

$$= \frac{\max_{i=1}^n \mathbf{x}_i + \frac{1}{\max_{i=2}^n \mathbf{x}_i} \frac{i}{1} \mathbf{x}_i}{\max_{i=1}^n \mathbf{x}_i + \frac{1}{\max_{i=2}^n \mathbf{x}_i} \frac{i}{1} \mathbf{x}_i} \frac{\mathbf{x}}{\max(\mathbf{x})}(k).$$

这说明规范化向量序列收敛到主特征值对应的特征向量.

同理, 可得到

$$\mathbf{v}_k = \frac{\max_{i=1}^{k-1} \mathbf{x}_i + \frac{1}{\max_{i=2}^n \mathbf{x}_i} \frac{i}{1} \mathbf{x}_i}{\max_{i=1}^{k-1} \mathbf{x}_i + \frac{1}{\max_{i=2}^n \mathbf{x}_i} \frac{i}{1} \mathbf{x}_i},$$

$$\max(\mathbf{v}_k) = \frac{\max_{i=1}^n \mathbf{x}_i + \frac{1}{\max_{i=2}^{k-1} \mathbf{x}_i} \frac{i}{1} \mathbf{x}_i}{\max_{i=1}^n \mathbf{x}_i + \frac{1}{\max_{i=2}^{k-1} \mathbf{x}_i} \frac{i}{1} \mathbf{x}_i}(k),$$

收敛速度由比值 $r = |\lambda_2/\lambda_1|$ 确定. 总结上述讨论, 有

定理 13 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 有 n 个线性无关的特征向量, 主特征值 λ_1 满足 $|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq \dots \geq |\lambda_n|$, 则对任意非零初始向量 $\mathbf{v} = \mathbf{u}_{(1)} \neq 0$, 按下述方法构造的向量序列 $\{\mathbf{u}\}, \{\mathbf{v}\}$:

$$\begin{aligned} \mathbf{v} &= \mathbf{u} - \mathbf{0}, \\ \mathbf{v}_k &= \mathbf{A}\mathbf{u}_{k-1}, \\ \mu_k &= \max(\mathbf{v}_k), \quad (k = 1, 2, \dots), \\ \mathbf{u} &= \mathbf{v}/\mu_k. \end{aligned} \tag{2.9}$$

则有

$$(1) \lim_k \mathbf{u}_k = \frac{\mathbf{x}}{\max(\mathbf{x})},$$

$$(2) \lim_k \mu_k = \lambda_1.$$

例 3 用幂法计算

$$\mathbf{A} = \begin{pmatrix} 1.0 & 1.0 & 0.5 \\ 1.0 & 1.0 & 0.25 \\ 0.5 & 0.25 & 2.0 \end{pmatrix}$$

的主特征值和相应的特征向量. 计算过程如表 8-1.

下述结果是用 8 位浮点数字进行运算得到的, \mathbf{u}_k 的分量值是舍入值. 于是得到

$$v_1 = 2.5365323$$

及相应的特征向量 $(0.7482, 0.6497, 1)^T$. 和相应的特征向量的真值(8 位数字)为

$$v_1 = 2.5365258,$$

$$\mathbf{u}_k = (0.74822116, 0.64966116, 1)^T.$$

表 8-1

k	\mathbf{u}_k^T (规范化向量)	$\max(\mathbf{v}_k)$
0	(1 1 1)	
1	(0.9091 0.8182 1)	2.7500000
5	(0.7651 0.6674 1)	2.5587918
10	(0.7494 0.6508 1)	2.5380029
15	(0.7483 0.6497 1)	2.5366256
16	(0.7483 0.6497 1)	2.5365840
17	(0.7482 0.6497 1)	2.5365598
18	(0.7482 0.6497 1)	2.5365456
19	(0.7482 0.6497 1)	2.5365374
20	(0.7482 0.6497 1)	2.5365323

8.2.2 加速方法

原点平移法

由前面讨论知道, 应用幂法计算 \mathbf{A} 的主特征值的收敛速度主

要由比值 $r = \frac{r_2}{r_1}$ 来决定, 但当 r 接近于 1 时, 收敛可能很慢. 这时, 一个补救的办法是采用加速收敛的方法.

引进矩阵

$$\mathbf{B} = \mathbf{A} - p\mathbf{I},$$

其中 p 为选择参数. 设 \mathbf{A} 的特征值为 $\lambda_1, \lambda_2, \dots, \lambda_n$, 则 \mathbf{B} 的相应特征值为 $\lambda_1 - p, \lambda_2 - p, \dots, \lambda_n - p$, 而且 \mathbf{A}, \mathbf{B} 的特征向量相同.

如果需要计算 \mathbf{A} 的主特征值 λ_1 , 就要适当选择 p 使 $\lambda_1 - p$ 仍然是 \mathbf{B} 的主特征值, 且使

$$\left| \frac{\lambda_2 - p}{\lambda_1 - p} \right| < \left| \frac{\lambda_2}{\lambda_1} \right|.$$

对 \mathbf{B} 应用幂法, 使得在计算 \mathbf{B} 的主特征值 $\lambda_1 - p$ 的过程中得到加速. 这种方法通常称为原点平移法. 对于 \mathbf{A} 的特征值的某种分布, 它是十分有效的.

例 4 设 $\mathbf{A} \in \mathbb{R}^{4 \times 4}$ 有特征值

$$\lambda_j = 15 - j \quad (j = 1, 2, 3, 4),$$

比值 $r = \frac{\lambda_2}{\lambda_1} = 0.9$. 作变换

$$\mathbf{B} = \mathbf{A} - p\mathbf{I} \quad (p = 12),$$

则 \mathbf{B} 的特征值为

$$\mu_1 = 2, \mu_2 = 1, \mu_3 = 0, \mu_4 = -1.$$

应用幂法计算 \mathbf{B} 的主特征值 μ_1 的收敛速度的比值为

$$\left| \frac{\mu_2}{\mu_1} \right| = \left| \frac{\lambda_2 - p}{\lambda_1 - p} \right| = \frac{1}{2} < \left| \frac{\lambda_2}{\lambda_1} \right| = 0.9.$$

虽然常常能够选择有利的 p 值, 使幂法得到加速, 但设计一个自动选择适当参数 p 的过程是困难的.

下面考虑当 \mathbf{A} 的特征值是实数时, 怎样选择 p 使采用幂法计

算₁ 得到加速 .

设 \mathbf{A} 的特征值满足

$$\lambda_1 > \lambda_2 > \dots > \lambda_n, \quad (2.10)$$

则不管 p 如何, $\mathbf{B} = \mathbf{A} - p\mathbf{I}$ 的主特征值为 $\lambda_1 - p$ 或 $\lambda_n - p$. 当我们希望计算 λ_1 及 \mathbf{x}_1 时, 首先应选择 p 使

$$|\lambda_1 - p| > |\lambda_n - p|,$$

且使收敛速度的比值

$$= \max \frac{|\lambda_2 - p|}{|\lambda_1 - p|}, \frac{|\lambda_n - p|}{|\lambda_1 - p|} = \min.$$

显然, 当 $\frac{\lambda_2 - p}{\lambda_1 - p} = -\frac{\lambda_n - p}{\lambda_1 - p}$, 即 $p = \frac{\lambda_2 + \lambda_n}{2}$ 时为最小, 这

时收敛速度的比值为

$$\frac{\lambda_2 - p^*}{\lambda_1 - p^*} = -\frac{\lambda_n - p^*}{\lambda_1 - p^*} = \frac{\lambda_2 - \lambda_n}{2\lambda_1 - \lambda_2 - \lambda_n}.$$

当 \mathbf{A} 的特征值满足 (2.10) 且 λ_2, λ_n 能初步估计时, 我们就能确定 p^* 的近似值 .

当希望计算 λ_n 时, 应选择

$$p = \frac{\lambda_1 + \lambda_{n-1}}{2} = p^*,$$

使得应用幂法计算 λ_n 得到加速 .

例 5 计算例 3 中矩阵 \mathbf{A} 的主特征值 .

作变换 $\mathbf{B} = \mathbf{A} - p\mathbf{I}$, 取 $p = 0.75$, 则

$$\begin{aligned} & 0.25 \quad 1 \quad 0.5 \\ \mathbf{B} = & \begin{pmatrix} 1 & 0.25 & 0.25 \\ 0.5 & 0.25 & 1.25 \end{pmatrix}. \end{aligned}$$

对 \mathbf{B} 应用幂法, 计算结果如表 8-2 .

表 8-2

k	\mathbf{u}_k^T (规范化向量)			$\max(\mathbf{v}_k)$
0	(1	1	1)	
5	(0 .7516	0 .6522	1)	1 .7914011
6	(0 .7491	0 .6511	1)	1 .7888443
7	(0 .7488	0 .6501	1)	1 .7873300
8	(0 .7484	0 .6499	1)	1 .7869152
9	(0 .7483	0 .6497	1)	1 .7866587
10	(0 .7482	0 .6497	1)	1 .7865914

由此得 \mathbf{B} 的主特征值为 $\mu_1 = 1 .7865914$, \mathbf{A} 的主特征值 λ_1 为

$$\lambda_1 = \mu_1 + 0 .75 = 2 .5365914,$$

与例 3 结果比较, 上述结果比例 3 迭代 15 次还好. 若迭代 15 次, $\mu = 1.7865258$ (相应的 $\lambda_1 = 2.5365258$) .

原点位移的加速方法, 是一个矩阵变换方法. 这种变换容易计算, 又不破坏矩阵 \mathbf{A} 的稀疏性, 但 p 的选择依赖于对 \mathbf{A} 的特征值分布的大致了解.

瑞利商加速

由定理 11 知, 对称矩阵 \mathbf{A} 的 λ_1 及 λ_n 可用瑞利商的极值来表示. 下面我们将把瑞利商应用到用幂法计算实对称矩阵 \mathbf{A} 的主特征值的加速收敛上来.

定理 14 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 为对称矩阵, 特征值满足

$$|\lambda_1| > |\lambda_2| > |\lambda_3| > \dots > |\lambda_n|,$$

对应的特征向量满足 $(\mathbf{x}_1, \mathbf{x}_2) = \delta_{ij}$, 应用幂法(公式(2.9))计算 \mathbf{A} 的主特征值 λ_1 , 则规范化向量 \mathbf{u}_k 的瑞利商给出 λ_1 的较好的近似

$$\frac{(\mathbf{A}\mathbf{u}_k, \mathbf{u}_k)}{(\mathbf{u}_k, \mathbf{u}_k)} = \lambda_1 + O\left(\frac{\lambda_2}{\lambda_1}\right)^{2k}.$$

证明 由(2.8)式及

$$\mathbf{u} = \frac{\mathbf{A}^k \mathbf{u}}{\max(\mathbf{A}^k \mathbf{u})}, \quad \mathbf{v}_{k+1} = \mathbf{A} \mathbf{u}_k = \frac{\mathbf{A}^{k+1} \mathbf{u}}{\max(\mathbf{A}^k \mathbf{u})},$$

得

$$\begin{aligned} \frac{(\mathbf{A} \mathbf{u}_k, \mathbf{u}_k)}{(\mathbf{u}, \mathbf{u}_k)} &= \frac{(\mathbf{A}^{k+1} \mathbf{u}, \mathbf{A}^k \mathbf{u})}{(\mathbf{A}^k \mathbf{u}, \mathbf{A}^k \mathbf{u})} = \frac{\sum_{j=1}^n \frac{2}{j} \frac{2}{j} \frac{k+1}{j}}{\sum_{j=1}^n \frac{2}{j} \frac{2}{j} \frac{k}{j}} \\ &= 1 + O\left(\frac{2}{1}\right)^{2k}. \end{aligned} \quad (2.11)$$

8.2.3 反幂法

反幂法用来计算矩阵按模最小的特征值及其特征向量,也可用来计算对应于一个给定近似特征值的特征向量.

设 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 为非奇异矩阵, \mathbf{A} 的特征值次序记为

$$\lambda_1 / \lambda_2 / \dots / \lambda_n /,$$

相应的特征向量为 $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, 则 \mathbf{A}^{-1} 的特征值为

$$\left| \frac{1}{\lambda_1} \right| \quad \left| \frac{1}{\lambda_2} \right| \quad \dots \quad \left| \frac{1}{\lambda_n} \right|,$$

对应的特征向量为 $\mathbf{x}_1, \mathbf{x}_{n-1}, \dots, \mathbf{x}_n$.

因此计算 \mathbf{A} 的按模最小的特征值 λ_n 的问题就是计算 \mathbf{A}^{-1} 的按模最大的特征值的问题.

对于 \mathbf{A}^{-1} 应用幂法迭代(称为反幂法), 可求得矩阵 \mathbf{A}^{-1} 的主特征值 λ_n , 从而求得 \mathbf{A} 的按模最小的特征值 λ_n .

反幂法迭代公式为:

任取初始向量 $\mathbf{v} = \mathbf{u} - \mathbf{0}$, 构造向量序列

$$\mathbf{v}_k = \mathbf{A}^{-1} \mathbf{u}_{k-1}$$

$$\mathbf{u}_k = \frac{\mathbf{v}_k}{\max(\mathbf{v}_k)} \quad (k = 1, 2, \dots).$$

迭代向量 \mathbf{v}_k 可以通过解方程组

$$\mathbf{A}\mathbf{v}_k = \mathbf{u}_{k-1}$$

求得.

定理 15 设 \mathbf{A} 为非奇异矩阵且有 n 个线性无关的特征向量, 其对应的特征值满足

$$|\lambda_1| > |\lambda_2| > \dots > |\lambda_{n-1}| > |\lambda_n| > 0,$$

则对任何初始非零向量 \mathbf{u}_0 ($\mathbf{u}_0 \neq 0$), 由反幂法构造的向量序列 $\{\mathbf{v}_k\}, \{\mathbf{u}_k\}$ 满足

$$(1) \lim_k \mathbf{u}_k = \frac{\mathbf{x}_n}{\max(\mathbf{x})},$$

$$(2) \lim_k \max(\mathbf{v}_k) = \frac{1}{\lambda_n}.$$

收敛速度的比值为 $\left| \frac{\lambda_n}{\lambda_{n-1}} \right|$.

在反幂法中也可以用原点平移法来加速迭代过程或求其他特征值及特征向量.

如果矩阵 $(\mathbf{A} - p\mathbf{I})^{-1}$ 存在, 显然其特征值为

$$\frac{1}{\lambda_1 - p}, \frac{1}{\lambda_2 - p}, \dots, \frac{1}{\lambda_n - p},$$

对应的特征向量仍然是 $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$. 现对矩阵 $(\mathbf{A} - p\mathbf{I})^{-1}$ 应用幂法, 得到反幂法的迭代公式

$$\mathbf{u} = \mathbf{v}_0 = \mathbf{0}, \text{ 初始向量}$$

$$\mathbf{v}_k = (\mathbf{A} - p\mathbf{I})^{-1} \mathbf{u}_{k-1} \quad (k = 1, 2, \dots). \quad (2.12)$$

$$\mathbf{u} = \frac{\mathbf{v}_k}{\max(\mathbf{v}_k)}$$

如果 p 是 \mathbf{A} 的特征值 λ_j 的一个近似值, 且设 λ_j 与其他特征值是分离的, 即

$$|\lambda_j - p| > |\lambda_i - p| \quad (i \neq j),$$

就是说 $\frac{1}{\lambda_j - p}$ 是 $(\mathbf{A} - p\mathbf{I})^{-1}$ 的主特征值, 可用反幂法(2.12)计算特

征值及特征向量 .

设 $A \in \mathbb{R}^{n \times n}$ 有 n 个线性无关的特征向量 $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, 则

$$\begin{aligned}\mathbf{u} &= \sum_{i=1}^n \lambda_i \mathbf{x}_i = (\lambda_1 \mathbf{x}_1 \quad \lambda_2 \mathbf{x}_2 \quad \dots \quad \lambda_n \mathbf{x}_n)^T, \\ \mathbf{v}_k &= \frac{(\mathbf{A} - p\mathbf{I})^{-k} \mathbf{u}}{\max((\mathbf{A} - p\mathbf{I})^{-k} \mathbf{u})}, \\ \mathbf{u}' &= \frac{(\mathbf{A} - p\mathbf{I})^{-k} \mathbf{u}}{\max((\mathbf{A} - p\mathbf{I})^{-k} \mathbf{u}')},\end{aligned}$$

其中

$$(\mathbf{A} - p\mathbf{I})^{-k} \mathbf{u} = \sum_{i=1}^n \lambda_i (\mathbf{x}_i - p\mathbf{x}_i)^T \mathbf{u}.$$

同理可得:

定理 16 设 $A \in \mathbb{R}^{n \times n}$ 有 n 个线性无关的特征向量, A 的特征值及对应的特征向量分别记为 λ_i 及 \mathbf{x}_i ($i = 1, 2, \dots, n$), 而 p 为 λ_j 的近似值, $(\mathbf{A} - p\mathbf{I})^{-1}$ 存在, 且

$$|\lambda_j - p| \approx |\lambda_i - p| \quad (i \neq j).$$

则对任意的非零初始向量 $\mathbf{u}_0 \neq 0$, 由反幂法迭代公式(2.12)构造的向量序列 $\{\mathbf{v}_k\}, \{\mathbf{u}_k\}$ 满足

$$(1) \lim_k \mathbf{u}_k = \frac{\mathbf{x}_j}{\max(\mathbf{x}_i)};$$

$$(2) \lim_k \max(\mathbf{v}_k) = \frac{1}{|\lambda_j - p|}, \text{ 即}$$

$$p + \frac{1}{\max(\mathbf{v}_k)} = \lambda_j \quad (\text{当 } k \rightarrow \infty),$$

且收敛速度由比值 $r = |\lambda_j - p| / \min_{i \neq j} |\lambda_i - p|$ 确定 .

由该定理知, 对 $\mathbf{A} - p\mathbf{I}$ (其中 $p \neq \lambda_j$) 应用反幂法, 可用来计算特征向量 \mathbf{x}_j . 只要选择的 p 是 λ_j 的一个较好的近似且特征值分离情况较好, 一般 r 很小, 常常只要迭代一二次就可完成特征向量的计算 .

反幂法迭代公式中的 \mathbf{v}_k 是通过解方程组

$$(\mathbf{A} - p\mathbf{I}) \mathbf{v}_k = \mathbf{u}_{k-1}$$

求得的.为了节省工作量,可以先将 $\mathbf{A} - p\mathbf{I}$ 进行三角分解

$$\mathbf{P}(\mathbf{A} - p\mathbf{I}) = \mathbf{L}\mathbf{U},$$

其中 \mathbf{P} 为某个排列阵,于是求 \mathbf{v}_k 相当于解两个三角形方程组

$$\mathbf{Ly}_k = \mathbf{Pu}_{k-1},$$

$$\mathbf{Uv}_k = \mathbf{y}_k.$$

实验表明,按上述方法选择 \mathbf{u} 是较好的:选 \mathbf{u} 使

$$\mathbf{Uv} = \mathbf{L}^{-1} \mathbf{Pu} = (1, 1, \dots, 1)^T \quad (2.13)$$

用回代求解(2.13)即得 \mathbf{v} ,然后再按公式(2.12)进行迭代.

反幂法计算公式

1. 分解计算

$\mathbf{P}(\mathbf{A} - p\mathbf{I}) = \mathbf{L}\mathbf{U}$,且保存 \mathbf{L}, \mathbf{U} 及 \mathbf{P} 信息.

2. 反幂法迭代

(1) 解 $\mathbf{Uv} = (1, \dots, 1)^T$ 求 \mathbf{v}

$$\mu_1 = \max(\mathbf{v}), \quad \mathbf{u} = \mathbf{v}/\mu_1$$

(2) $k = 2, 3, \dots$

1) 解 $\mathbf{Ly}_k = \mathbf{Pu}_{k-1}$ 求 \mathbf{y}_k

解 $\mathbf{Uv}_k = \mathbf{y}_k$ 求 \mathbf{v}_k

2) $\mu_k = \max(\mathbf{v}_k)$

3) 计算 $\mathbf{u}_k = \mathbf{v}_k/\mu_k$

例 6 用反幂法求

$$\begin{aligned} & 2 & 1 & 0 \\ \mathbf{A} = & 1 & 3 & 1 \\ & 0 & 1 & 4 \end{aligned}$$

的对应于计算特征值 $\lambda = 1.2679$ (精确特征值为 $\lambda_3 = 3 - \sqrt{3}$) 的特征向量(用 5 位浮点数进行运算).

解 用部分选主元的三角分解将 $\mathbf{A} - p\mathbf{I}$ (其中 $p = 1.2679$) 分

解为

$$\mathbf{P}(\mathbf{A} - p\mathbf{I}) = \mathbf{LU},$$

其中

$$\begin{aligned}\mathbf{L} &= \begin{matrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0.7321 & -0.26807 & 1 \end{matrix}, \\ \mathbf{U} &= \begin{matrix} 1 & 1.7321 & 1 \\ 0 & 1 & 2.7321 \\ 0 & 0 & 0.29405 \times 10^{-3} \end{matrix}, \\ \mathbf{P} &= \begin{matrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{matrix}.\end{aligned}$$

由 $\mathbf{Uv} = (1, 1, 1)^T$, 得

$$\begin{aligned}\mathbf{v} &= (12692, -9290.3, 3400.8)^T, \\ \mathbf{u} &= (1, -0.73198, 0.26795)^T,\end{aligned}$$

由 $\mathbf{LUv} = \mathbf{Pu}$, 得

$$\begin{aligned}\mathbf{v} &= (20404, -14937, 5467.4)^T, \\ \mathbf{u} &= (1, -0.73206, 0.26796)^T,\end{aligned}$$

\mathbf{x}_3 对应的特征向量是

$$\mathbf{x} = (1, 1 - 3, 2 - 3)^T = (1, -0.73205, 0.26795)^T,$$

由此看出 \mathbf{u} 是 \mathbf{x} 的相当好的近似.

特征值 $\lambda_3 = 1.2679 + 1/\mu = 1.26794901$, \mathbf{x}_3 的真值为 $\mathbf{x}_3 = 3 - 3 = 1.26794912\dots$.

8.3 豪斯霍尔德方法

8.3.1 引言

本节讨论两个问题

(1) 用初等反射阵作正交相似变换约化一般实矩阵 \mathbf{A} 为上海森伯格阵 .

(2) 用初等反射阵作正交相似变换约化对称矩阵 \mathbf{A} 为对称三对角阵 .

于是, 求原矩阵特征值问题, 就转化为求上海森伯格阵或对称三对角阵的特征值问题 .

8.3.2 用正交相似变换约化一般矩阵为上海森伯格阵

设 $\mathbf{A} = (a_{ij}) \in \mathbf{R}^{n \times n}$. 下面来说明, 可选择初等反射阵 $\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_{n-1}$ 使 \mathbf{A} 经正交相似变换约化为一个上海森伯格阵 .

(1) 设

$$\mathbf{A} = \begin{matrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \cdots & \cdots & \cdots & \cdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{matrix} = \begin{matrix} a_{11} & \mathbf{A}_{12}^{(1)} \\ \mathbf{c} & \mathbf{A}_{22}^{(1)} \end{matrix},$$

其中 $\mathbf{c} = (a_{21}, \dots, a_{n1})^T \in \mathbf{R}^{n-1}$, 不妨设 $\mathbf{c} \neq \mathbf{0}$, 否则这一步不需要约化. 于是, 可选择初等反射阵 $\mathbf{R} = \mathbf{I} - \frac{1}{\| \mathbf{c} \|_2^2} \mathbf{c} \mathbf{c}^T$ 使 $\mathbf{R} \mathbf{c} = -e_1$, 其中

$$\begin{aligned} e_1 &= \operatorname{sgn}(a_{21}) \sqrt{\sum_{i=2}^n a_{ii}^2} e_1^{1/2}, \\ \mathbf{u}_1 &= \mathbf{c} + e_1, \\ e_1 &= -e_1(e_1 + a_{21}). \end{aligned} \tag{3.1}$$

令

$$\mathbf{U}_1 = \frac{1}{\|\mathbf{c}\|_2},$$

则

$$\mathbf{A}_1 = \mathbf{U}_1 \mathbf{A} \mathbf{U}_1 = \frac{a_{11} \quad \mathbf{A}_{12}^{(1)} \mathbf{R}}{\mathbf{R} \mathbf{c} \quad \mathbf{R} \mathbf{A}_{22}^{(1)} \mathbf{R}}$$

$$\begin{array}{cccccc}
 & a_{11} & a_{12}^{(2)} & a_{13}^{(2)} & \dots & a_{1n}^{(2)} \\
 - & 1 & a_{22}^{(2)} & a_{23}^{(2)} & \dots & a_{2n}^{(2)} \\
 = & 0 & a_{32}^{(2)} & a_{33}^{(2)} & \dots & a_{3n}^{(2)} & \left| \begin{array}{c} \mathbf{A}_{11}^{(2)} \\ \mathbf{A}_{22}^{(2)} \\ \mathbf{A}_{32}^{(2)} \end{array} \right. \\
 & \dots & \dots & \dots & \dots & \dots \\
 & 0 & a_{n2}^{(2)} & a_{n3}^{(2)} & \dots & a_{nn}^{(2)}
 \end{array} ,$$

其中 $\mathbf{c} = (a_{12}^{(2)}, \dots, a_{n2}^{(2)})^T \in \mathbb{R}^{n-2}$, $\mathbf{A}_{22}^{(2)} \in \mathbb{R}^{(n-2) \times (n-2)}$.

(2) 第 k 步约化: 重复上述过程, 设对 \mathbf{A} 已完成第 1 步, ..., 第 $k-1$ 步正交相似变换, 即有

$$\mathbf{A}_k = \mathbf{U}_{k-1} \mathbf{A}_{k-1} \mathbf{U}_{k-1},$$

或

$$\mathbf{A}_k = \mathbf{U}_{k-1} \dots \mathbf{U}_1 \mathbf{A}_1 \mathbf{U}_1 \dots \mathbf{U}_{k-1},$$

且

$$\begin{array}{ccccccccc}
 & a_{11}^{(1)} & a_{12}^{(2)} & \dots & a_{1,k-1}^{(k-1)} & a_{1k}^{(k)} & a_{1,k+1}^{(k)} & \dots & a_{1n}^{(k)} \\
 - & 1 & a_{22}^{(2)} & \dots & a_{2,k-1}^{(k-1)} & a_{2k}^{(k)} & a_{2,k+1}^{(k)} & \dots & a_{2n}^{(k)} \\
 & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\
 \mathbf{A}_k = & & & & a_{kk}^{(k)} & a_{k,k+1}^{(k)} & \dots & a_{kn}^{(k)} \\
 & & & & a_{k+1,k}^{(k)} & a_{k+1,k+1}^{(k)} & \dots & a_{k+1,n}^{(k)} \\
 & & & & \vdots & \vdots & \ddots & \vdots \\
 & & & & a_{n,k}^{(k)} & a_{n,k+1}^{(k)} & \dots & a_{nn}^{(k)} \\
 & k & n-k & & & & & & \\
 & \mathbf{A}_1^{(k)} & \mathbf{A}_2^{(k)} & \dots & \mathbf{A}_k^{(k)} & & & & \\
 & \mathbf{0} \quad \mathbf{c} & \mathbf{A}_{22}^{(k)} & n-k, & & & & &
 \end{array}$$

其中 $\mathbf{c} = (a_{k+1,k}^{(k)}, \dots, a_{n,k}^{(k)})^T \in \mathbb{R}^{n-k}$, $\mathbf{A}_1^{(k)}$ 为 k 阶上三角矩阵, $\mathbf{A}_2^{(k)}$ 为 $(n-k) \times (n-k)$.

设 $\mathbf{c} \neq \mathbf{0}$, 于是可选择初等反射阵 \mathbf{R}_k 使 $\mathbf{R}_k \mathbf{c} = -e_k$, 其中, \mathbf{R}_k 计算公式为

$$k = \operatorname{sgn}(a_{k+1,k}^{(k)}) \prod_{i=k+1}^n (a_{ik}^{(k)})^2^{-1/2},$$

$$\begin{aligned}\mathbf{u}_k &= \mathbf{c}_k + e_k \mathbf{e}_k, \\ k &= a_{k+1,k}^{(k)} + e_k, \\ \mathbf{R}_k &= \mathbf{I} - e_k \mathbf{u}_k \mathbf{u}_k^T.\end{aligned}\quad (3.2)$$

令

$$\mathbf{U}_k = \frac{\mathbf{I}}{\mathbf{R}_k},$$

则

$$\begin{aligned}\mathbf{A}_{k+1} &= \mathbf{U}_k \mathbf{A}_k \mathbf{U}_k = \left| \begin{array}{c|cc} \mathbf{A}_1^{(k+1)} & \mathbf{A}_{12}^{(k)} & \mathbf{R}_k \\ \mathbf{0} & \mathbf{R}_k \mathbf{c}_k & \mathbf{R}_k \end{array} \right| \mathbf{R}_k \mathbf{A}_{22}^{(k)} \mathbf{R}_k \\ &= \left| \begin{array}{c|cc} \mathbf{A}_{11}^{(k+1)} & \mathbf{A}_{12}^{(k+1)} & \\ \mathbf{0} & \mathbf{c}_{k+1} & \mathbf{A}_{22}^{(k+1)} \end{array} \right|, \quad (3.3)\end{aligned}$$

其中 $\mathbf{A}_1^{(k+1)}$ 为 $k+1$ 阶上海森伯格阵. 第 k 步约化只需计算 $\mathbf{A}_1^{(k)}$ \mathbf{R}_k 及 $\mathbf{R}_k \mathbf{A}_{22}^{(k)} \mathbf{R}_k$ (当 \mathbf{A} 为对称阵时, 只需计算 $\mathbf{R}_k \mathbf{A}_{22}^{(k)} \mathbf{R}_k$).

(3) 重复上述过程, 则有

$$\begin{aligned}&\mathbf{U}_{n-2} \dots \mathbf{U}_2 \mathbf{U}_1 \mathbf{A} \mathbf{U}_1 \mathbf{U}_2 \dots \mathbf{U}_{n-2} \\&\begin{array}{ccccccc} a_{11} & * & * & \dots & * & * \\ -1 & \mathbf{d}_2^{(2)} & * & \dots & * & * \\ -2 & \mathbf{d}_3^{(3)} & * & \dots & * & * \\ \vdots & & & & & \\ -n-2 & & & & & & * \\ -n-1 & & & & & & \mathbf{d}_{n-1,n-1}^{(n-2)} \\ & & & & & & \\ & & & & & & \mathbf{d}_{nn}^{(n-1)} \end{array} \\&= \mathbf{A}_{n-1}.\end{aligned}$$

总结上述讨论, 有

定理 17 (豪斯霍尔德约化矩阵为上海森伯格阵) 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$, 则存在初等反射阵 $\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_{n-2}$ 使

$$\mathbf{U}_{n-2} \dots \mathbf{U}_2 \mathbf{U}_1 \mathbf{A} \mathbf{U}_1 \mathbf{U}_2 \dots \mathbf{U}_{n-2} \quad \mathbf{U}_0^T \mathbf{A} \mathbf{U}_0 = \mathbf{H}(\text{上海森伯格阵}).$$

算法 1 (豪斯霍尔德约化矩阵为上海森伯格型) 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$,

本算法计算 $\mathbf{U}^T \mathbf{A} \mathbf{U} = \mathbf{H}$ (上海森伯格型), 其中 $\mathbf{U}_0 = \mathbf{U}_1 \mathbf{U}_2 \dots \mathbf{U}_{n-2}$ 为初等反射阵的乘积 .

1. $\mathbf{U}_0 = \mathbf{I}$

2. 对于 $k=1, 2, \dots, n-2$

(1) 计算初等反射阵 \mathbf{R}_k 使

$$\mathbf{R}_k \mathbf{c}_k = -\mathbf{e}_k$$

(2) 约化计算

$$\mathbf{A} = \mathbf{U}_k \mathbf{A} \mathbf{U}_k^T, \quad \mathbf{U}_k = \begin{matrix} & & & \mathbf{I}_k \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \mathbf{R}_k \end{matrix}$$

(3) $\mathbf{U} = \mathbf{U}_0 \mathbf{U}_k$

本算法约需要 $\frac{5}{3} n^3$ 次乘法运算, 要明显形成 \mathbf{U}_0 还需要附加

$\frac{2}{3} n^3$ 次乘法 .

例 7 用豪斯霍尔德方法将

$$\mathbf{A} = \begin{matrix} & -4 & -3 & -7 \\ & 2 & 3 & 2 \\ & 4 & 2 & 7 \end{matrix}$$

矩阵约化为上 Hessenberg 阵 .

解 选取初等反射阵 \mathbf{R} 使 $\mathbf{R} \mathbf{c} = -\mathbf{e}_1$, 其中 $\mathbf{c} = (2, 4)^T$.

(1) 计算 \mathbf{R} : $\|\mathbf{c}\|_2 = \max(2, 4) = 4$, $\mathbf{c}' = \mathbf{c} / \|\mathbf{c}\|_2 = (0.5, 1)^T$ (规范化)

$$= 1.25 = 1.118034,$$

$$\mathbf{u} = \mathbf{c} + \mathbf{e} = (1.118034, 1)^T,$$

$$\mathbf{r}_1 = (\mathbf{u} + 0.5) = 1.809017,$$

$$\mathbf{r}_2 = \mathbf{u} = 4.472136,$$

$$\mathbf{R} = \mathbf{I} - \mathbf{r}_1^{-1} \mathbf{u} \mathbf{u}^T.$$

则有

$$\mathbf{R} \mathbf{c} = -\mathbf{e}_1.$$

(2) 约化计算:

令

$$\mathbf{U}_1 = \begin{pmatrix} 1 & \mathbf{0} \\ \mathbf{0} & \mathbf{R} \end{pmatrix},$$

则

$$\mathbf{A} = \mathbf{U}_1 \mathbf{A} \mathbf{U}_1 = \begin{pmatrix} -4 & & & \\ & 7.602631 & -0.447214 & \\ -4.472136 & & 7.799999 & -0.400000 \\ 0 & & -0.399999 & 2.200000 \end{pmatrix} = \mathbf{H}.$$

8.3.3 用正交相似变换约化对称阵为对称三对角阵

定理 18 (豪斯霍尔德约化对称阵为对称三对角阵) 设
A $\in \mathbb{R}^{n \times n}$ 为对称矩阵, 则存在初等反射阵 $\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_{n-2}$ 使

$$\begin{array}{ccccc} a & & b & & \\ & b & c & b & \\ \mathbf{U}_{n-2} \dots \mathbf{U}_2 \mathbf{U}_1 \mathbf{A} \mathbf{U}_1 \mathbf{U}_2 \dots \mathbf{U}_{n-2} & = & w & w & w \\ & & & & \\ & & b_{n-2} & c_{n-1} & b_{n-1} \\ & & & & \\ & & b_{n-1} & c_n & \end{array} \quad \mathbf{C}.$$

证明 由定理 17, 存在初等反射阵 $\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_{n-2}$ 使 $\mathbf{U}_{n-2} \dots \mathbf{U}_2 \mathbf{U}_1 \mathbf{A} \mathbf{U}_1 \mathbf{U}_2 \dots \mathbf{U}_{n-2} = \mathbf{H} = \mathbf{A}_{n-1}$ 为上三角阵, 且 \mathbf{A}_{n-1} 亦是对称阵, 因此, \mathbf{A}_{n-1} 为对称三对角阵.

由上面讨论可知, 当 \mathbf{A} 为对称阵时, 由 $\mathbf{A}_k \mathbf{A}_{k+1} = \mathbf{U}_k \mathbf{A}_k \mathbf{U}_k$ 一步约化计算中只需计算 \mathbf{R}_k 及 $\mathbf{R}_k \mathbf{A}_{22}^{(k)} \mathbf{R}_k$. 又由于 \mathbf{A} 的对称性, 故只需计算 $\mathbf{R}_k \mathbf{A}_{22}^{(k)} \mathbf{R}_k$ 的对角线以下元素. 注意到

$$\mathbf{R}_k \mathbf{A}_{22}^{(k)} \mathbf{R}_k = (\mathbf{I} - \frac{1}{k} \mathbf{u}_k \mathbf{u}_k^T) (\mathbf{A}_{22}^{(k)} - \frac{1}{k} \mathbf{A}_{22}^{(k)} \mathbf{u}_k \mathbf{u}_k^T).$$

引进记号

$$\mathbf{r}_k = \frac{1}{k} \mathbf{A}_{22}^{(k)} \mathbf{u}_k \in \mathbb{R}^{n-k},$$

$$\mathbf{t}_k = \mathbf{r}_k - \frac{1}{2} (\mathbf{u}_k^T \mathbf{r}_k) \mathbf{u}_k \in \mathbb{R}^{n-k},$$

则

$$\begin{aligned} \mathbf{R}_k \mathbf{A}_{22}^{(k)} \mathbf{R}_k &= \mathbf{A}_{22}^{(k)} - \mathbf{u}_k \mathbf{t}_k^T - \mathbf{t}_k \mathbf{u}_k^T \\ (i = k+1, \dots, n, j = k+1, \dots, i). \end{aligned}$$

算法 2(豪斯霍尔德约化对称阵为对称三对角阵) 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 为对称阵, 本算法确定初等反射阵 \mathbf{U}_k ($k = 1, \dots, n-2$) 使 $\mathbf{U}_{n-2} \dots \mathbf{U}_1 \mathbf{A} \mathbf{U}_{n-2} \dots \mathbf{U}_1 = \mathbf{C}$ (为对称三对角阵). \mathbf{C} 的对角元 c_i ($i = 1, \dots, n$) 存放在数组 $c(n)$ 内, \mathbf{C} 的次对角元素 b_i ($i = 1, \dots, n-1$) 存放在数组 $b(n)$ 内. 数组 $b(n)$ 最初可用来存放 \mathbf{r}_k 及 \mathbf{t}_k , 确定 \mathbf{R}_k 中向量 \mathbf{u} 的分量存放在 \mathbf{A} 的相应位置. a_{kk} 冲掉 a_{kk} , 约化 \mathbf{A} 的结果冲掉 \mathbf{A}_k , 数组 \mathbf{A} 的上部分元素不变. 如果第 k 步不需要变换则置 a_{kk} 为零.

1. 对于 $k = 1, 2, \dots, n-2$

(1) $c_k = a_{kk}$

(2) 确定变换 \mathbf{R}_k

1) 计算 $d = \max_{i=k+1}^n |a_{ik}|$

2) 如果 $d = 0$, 则 $a_{kk} = 0, b_k = 0$, 转 4)

3) 计算 $a_{ik} - u_{ik} = a_{ik}/d$ ($i = k+1, \dots, n$)

4) $u_{k+1, k} = \operatorname{sgn}(a_{k+1, k}) \sum_{i=k+1}^n u_{ik}^2$

5) $u_{k+1, k} = u_{k+1, k} +$

6) $a_{kk} - u_{kk} = -u_{k+1, k}^*$

7) $b_k - u_{kk} = -u_{k+1, k}^* d$

(3) 应用变换

1) $s = 0$

2) 计算 $\mathbf{A}_2^{(k)} \mathbf{u}$ 及 $\mathbf{u}_k^T \mathbf{r}_k$

对于 $i = k+1, \dots, n$

(a) $b_i - h = \sum_{j=k+1}^i a_{ij} * u_{jk} + \sum_{j=i+1}^n a_{ji} * u_{jk}$

(b) $s = s + h * u_{ik}$

3) 计算 t_k

$$b_k = (b_k - (s * u_{ik}) / (2 * u_{kk})) / u_{kk}$$

$$(i = k+1, \dots, n)$$

4) 计算 $\mathbf{R}_k \mathbf{A}_{22}^{(k)} \mathbf{R}_k$ 对角线以下部分

对于 $i = k+1, \dots, n, j = k+1, \dots, i$

$$a_{ij} = a_{ij} - u_{ik} * b_j - b_i * u_{jk}$$

5) 继续循环 k

2. $c_{n-1} = a_{n-1, n-1}$

$$c_n = a_{n, n}$$

$$b_{n-1} = a_{n, n-1}$$

对对称阵 \mathbf{A} 用初等反射阵正交相似约化为对称三对角阵大约需要 $\frac{2}{3} n^3$ 次乘法 .

用正交矩阵进行相似约化有一些特点, 如构造的 \mathbf{U}_k 容易求逆, 且 \mathbf{U}_k 的元素数量级不大, 这个算法是十分稳定的 .

8.4 QR 方法

8.4.1 QR 算法

Rutishauser(1958)利用矩阵的三角分解提出了计算矩阵特征值的 LR 算法, Francis(1961, 1962)利用矩阵的 QR 分解建立了计算矩阵特征值的 QR 方法 .

QR 方法是一种变换方法, 是计算一般矩阵(中小型矩阵)全部特征值问题的最有效方法之一 .

目前 QR 方法主要用来计算:(1)上海森伯格阵的全部特征值问题,(2)计算对称三对角矩阵的全部特征值问题, 且 QR 方法具有收敛快, 算法稳定等特点 .

对于一般矩阵 $\mathbf{A} \in \mathbb{R}^{n \times n}$ (或对称矩阵), 则首先用豪斯霍尔德方法将 \mathbf{A} 化为上海森伯格阵 \mathbf{B} (或对称三对角阵), 然后再用 QR

方法计算 \mathbf{B} 的全部特征值 .

设 $\mathbf{A} \in \mathbb{R}^{n \times n}$, 且对 \mathbf{A} 进行 QR 分解, 即

$$\mathbf{A} = \mathbf{Q}\mathbf{R},$$

其中 \mathbf{R} 为上三角阵, \mathbf{Q} 为正交阵, 于是可得到一新矩阵

$$\mathbf{B} = \mathbf{R}\mathbf{Q} = \mathbf{Q}^T \mathbf{A} \mathbf{Q}.$$

显然, \mathbf{B} 是由 \mathbf{A} 经过正交相似变换得到, 因此 \mathbf{B} 与 \mathbf{A} 特征值相同 . 再对 \mathbf{B} 进行 QR 分解, 又可得一新的矩阵, 重复这一过程可得到矩阵序列:

设 $\mathbf{A}_0 = \mathbf{A}$

将 \mathbf{A}_0 进行 QR 分解 $\mathbf{A}_0 = \mathbf{Q}_0 \mathbf{R}_0$

作矩阵 $\mathbf{A}_1 = \mathbf{R}_0 \mathbf{Q}_0 = \mathbf{Q}_0^T \mathbf{A}_0 \mathbf{Q}_0$

...

求得 \mathbf{A}_k 后将 \mathbf{A}_k 进行 QR 分解 $\mathbf{A}_k = \mathbf{Q}_k \mathbf{R}_k$

形成矩阵 $\mathbf{A}_{k+1} = \mathbf{R}_k \mathbf{Q}_k = \mathbf{Q}_k^T \mathbf{A}_k \mathbf{Q}_k$

...

QR 算法, 就是利用矩阵的 QR 分解, 按上述递推法则构造矩阵序列 $\{\mathbf{A}_k\}$ 的过程 . 只要 \mathbf{A} 为非奇异矩阵, 则由 QR 算法就完全确定 $\{\mathbf{A}_k\}$.

定理 19(基本 QR 方法) 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$. 构造 QR 算法:

$$\begin{aligned} \mathbf{A}_k &= \mathbf{Q}_k \mathbf{R}_k \quad (\text{其中 } \mathbf{Q}_k^T \mathbf{Q}_k = \mathbf{I}, \mathbf{R}_k \text{ 为上三角阵}); \\ \mathbf{A}_{k+1} &= \mathbf{R}_k \mathbf{Q}_k \quad (k = 1, 2, \dots), \end{aligned} \tag{4.1}$$

记 $\mathbf{Q} = \mathbf{Q}_1 \mathbf{Q}_2 \dots \mathbf{Q}_k$, $\mathbf{R} = \mathbf{R}_1 \mathbf{R}_2 \dots \mathbf{R}_k$, 则有

(1) \mathbf{A}_{k+1} 相似于 \mathbf{A}_k , 即 $\mathbf{A}_{k+1} = \mathbf{A}_k^T \mathbf{A}_k \mathbf{Q}_k$;

(2) $\mathbf{A}_{k+1} = (\mathbf{Q}_1 \mathbf{Q}_2 \dots \mathbf{Q}_k)^T \mathbf{A} (\mathbf{Q}_1 \mathbf{Q}_2 \dots \mathbf{Q}_k)$
 $= \mathbf{Q}^T \mathbf{A} \mathbf{Q}$;

(3) \mathbf{A}^k 的 QR 分解式为 $\mathbf{A}^k = \mathbf{Q}^k \mathbf{R}^k$.

证明 (1), (2) 显然, 现证 (3). 用归纳法, 显然, 当 $k = 1$ 时有

$\mathbf{A} = \mathbf{Q} \mathbf{R} = \mathbf{Q}_k \mathbf{R}_k$, 设 \mathbf{A}^{k-1} 有分解式

$$\mathbf{A}^{k-1} = \mathbf{Q}_{k-1} \mathbf{R}_{k-1},$$

于是

$$\begin{aligned}\mathbf{A} \mathbf{R}_k &= \mathbf{Q}_k \mathbf{Q}_{k-1} \dots (\mathbf{Q}_2 \mathbf{R}_2) \dots \mathbf{R}_k \\ &= \mathbf{Q}_k \mathbf{Q}_{k-1} \dots \mathbf{Q}_2 \mathbf{A}_k \mathbf{R}_{k-1} \dots \mathbf{R}_k \\ &= \mathbf{Q}_{k-1} \mathbf{A}_k \mathbf{R}_{k-1} \\ &= \mathbf{A}_k \mathbf{Q}_{k-1} \mathbf{R}_{k-1} = \mathbf{A}^k \text{ (因为 } \mathbf{A}_k = \mathbf{Q}_{k-1} \mathbf{A} \mathbf{Q}_{k-1}^T\text{)}.\end{aligned}$$

由第 5 章定理 30 或定理 31 知, 将 \mathbf{A}_k 进行 QR 分解, 即将 \mathbf{A}_k 用正交变换(左变换)化为上三角矩阵

$$\mathbf{Q}_k^T \mathbf{A}_k = \mathbf{R}_k,$$

其中 $\mathbf{Q}^T = \mathbf{P}_{n-1} \dots \mathbf{P}_2 \mathbf{P}_1$, 故

$$\mathbf{A}_{k+1} = \mathbf{Q}_k^T \mathbf{A}_k \mathbf{Q}_k = \mathbf{P}_{n-1} \dots \mathbf{P}_2 \mathbf{P}_1 \mathbf{A}_k \mathbf{P}_1^T \mathbf{P}_2^T \dots \mathbf{P}_{n-1}^T.$$

这就是说 \mathbf{A}_{k+1} 可由 \mathbf{A}_k 按下述方法求得:

(1) 左变换 $\mathbf{P}_{n-1} \dots \mathbf{P}_2 \mathbf{P}_1 \mathbf{A}_k = \mathbf{R}_k$ (上三角阵);

(2) 右变换 $\mathbf{R}_k \mathbf{P}_1^T \mathbf{P}_2^T \dots \mathbf{P}_{n-1}^T = \mathbf{A}_{k+1}$.

定理 20(QR 方法的收敛性) 设 $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{n \times n}$,

(1) 如果 \mathbf{A} 的特征值满足: $| \lambda_1 | > | \lambda_2 | > \dots > | \lambda_n | > 0$;

(2) \mathbf{A} 有标准型 $\mathbf{A} = \mathbf{X} \mathbf{D} \mathbf{X}^{-1}$ 其中 $\mathbf{D} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$, 且设 \mathbf{X}^{-1} 有三角分解 $\mathbf{X}^{-1} = \mathbf{L} \mathbf{U}$ (\mathbf{L} 为单位下三角阵, \mathbf{U} 为上三角阵), 则由 QR 算法产生的 $\{\mathbf{A}_k\}$ 本质上收敛于上三角矩阵, 即

$$\mathbf{A}_k \xrightarrow{\text{本质上}} \mathbf{R} = \begin{matrix} & * & & * \\ 1 & & \dots & \\ & 2 & \dots & * \\ & \vdots & \dots & \\ & n & & \end{matrix} \quad (\text{当 } k \rightarrow \infty \text{ 时})$$

若记 $\mathbf{A}_k = [a_{ij}^{(k)}]$, 则

$$(1) \lim_{k \rightarrow \infty} a_{ii}^{(k)} = \lambda_i; \quad (4.2)$$

(2) 当 $i > j$ 时, $\lim_k a_{ij}^{(k)} = 0$; (4.3)

当 $i < j$ 时 $a_{ij}^{(k)}$ 极限不一定存在.

证明可参阅 [1].

定理 21 如果对称矩阵 \mathbf{A} 满足定理 20 的条件, 则由 QR 算法产生的 $\{\mathbf{A}_k\}$ 收敛于对角阵 $\mathbf{D} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$.

证明 由定理 20 即知.

关于 QR 算法收敛性的进一步结果为:

设 $\mathbf{A} \in \mathbb{R}^{n \times n}$, 且 \mathbf{A} 有完备的特征向量集合, 如果 \mathbf{A} 的等模特征值中只有实重特征值或多重复的共轭特征值, 则由 QR 算法产生的 $\{\mathbf{A}_k\}$ 本质收敛于分块上三角矩阵 (对角块为一阶和二阶子块) 且对角块中每一个 2×2 子块给出 \mathbf{A} 的一对共轭复特征值, 每一个一阶对角子块给出 \mathbf{A} 的实特征值, 即

$$\begin{array}{ccccccccc} & & * & * & * & * & * & * & * \\ 1 & \dots \\ & & W & \dots & \dots & \dots & \dots & \dots & \dots \\ & & m & * & * & \dots & * & * & \\ \mathbf{A}_k & & & & & & & & \\ & & & & \mathbf{B} & & & & \\ & & & & & \dots & * & * & , \\ & & & & & & & & \\ & & & & & & W & \dots & \\ & & & & & & & & \\ & & & & & & & \mathbf{B} & \\ \end{array}$$

其中 $m + 2l = n$, \mathbf{B} 为 2×2 子块, 它给出 \mathbf{A} 一对共轭特征值.

8.4.2 带原点位移的 QR 方法

经分析指出: 定理 20 中 $\lim_k a_{nn}^{(k)} = \lambda_n$ 的速度依赖于比值 $r_n = |n' - n_{-1}|$, 当 r_n 很小时, 收敛较快, 如果 s 为 λ_n 的一个估计, 且对 $\mathbf{A} - s\mathbf{I}$ 运用 QR 算法, 则 $(n, n-1)$ 元素将以收敛因子

$|(-n - s)/(-n - 1 - s)|$ 线性收敛于零, (n, n) 元素将比在基本算法中收敛更快.

为此,为了加速收敛,选择数列 $\{s_k\}$,按上述方法构造矩阵序列 $\{\mathbf{A}_k\}$,称为带原点位移的QR算法.

设 $\mathbf{A} = \mathbf{A} \in \mathbf{R}^{n \times n}$

对 $\mathbf{A} - s_1 \mathbf{I}$ 进行 QR 分解 $\mathbf{A} - s_1 \mathbf{I} = \mathbf{Q} \mathbf{R}$

$$\begin{aligned}\text{形成矩阵 } \mathbf{A} &= \mathbf{R} \mathbf{Q} + s_1 \mathbf{I} = \mathbf{Q}^T (\mathbf{A} - s_1 \mathbf{I}) \mathbf{Q} + s_1 \mathbf{I} \\ &= \mathbf{Q}^T \mathbf{A} \mathbf{Q}\end{aligned}$$

求得 \mathbf{A}_k 后,将 $\mathbf{A}_k - s_k \mathbf{I}$ 进行 QR 分解

$$\mathbf{A}_k - s_k \mathbf{I} = \mathbf{Q}_k \mathbf{R}_k, \quad k = 3, 4, \dots \quad (4.4)$$

形成矩阵

$$\mathbf{A}_{k+1} = \mathbf{R}_k \mathbf{Q}_k + s_k \mathbf{I} = \mathbf{Q}_k^T \mathbf{A}_k \mathbf{Q}_k \quad (4.5)$$

如果令 $\mathbf{D}_k = \mathbf{Q}_k \mathbf{Q}_k^T \dots \mathbf{Q}_1 \mathbf{Q}_1^T$, $\mathbf{D}_k = \mathbf{R}_k \dots \mathbf{R}_1 \mathbf{R}_1^T$, 则有 $\mathbf{A}_{k+1} = \mathbf{D}_k^T \mathbf{A} \mathbf{D}_k$, 并且矩阵 $(\mathbf{A} - s_1 \mathbf{I})(\mathbf{A} - s_2 \mathbf{I}) \dots (\mathbf{A} - s_n \mathbf{I})$ (A) 有 QR 分解式

$$(A) = \mathbf{D}_n \mathbf{D}_n^T.$$

在带位移 QR 方法中,每步并不需要形成 \mathbf{Q} 和 \mathbf{R} ,可按下面的方法计算:

首先用正交变换(左变换)将 $\mathbf{A} - s_k \mathbf{I}$ 化为上三角阵,即

$$\mathbf{P}_{n-1} \dots \mathbf{P}_2 \mathbf{P}_1 (\mathbf{A} - s_k \mathbf{I}) = \mathbf{R}$$

(当 \mathbf{A} 为上海森伯格阵或对称三对角阵时, \mathbf{P}_i 可为平面旋转阵),则

$$\mathbf{A}_{k+1} = \mathbf{P}_{n-1} \dots \mathbf{P}_2 \mathbf{P}_1 (\mathbf{A} - s_k \mathbf{I}) \mathbf{P}_1^T \mathbf{P}_2^T \dots \mathbf{P}_{n-1}^T + s_k \mathbf{I}.$$

下面考虑用 QR 方法计算上海森伯格阵的特征值.

设 \mathbf{B} 为上海森伯格阵,即

$$\mathbf{B} = \begin{array}{cccc} b_1 & b_2 & \dots & b_n \\ & b_{11} & b_{22} & \dots & b_{nn} \\ & w & w & \dots & \\ & b_{n,n-1} & b_{n,n} & & \end{array}.$$

如果 $b_{i+1,i} = 0$ ($i = 1, 2, \dots, n-1$), 则称 \mathbf{B} 为不可约上海森伯格阵.

设 $\mathbf{A} \in \mathbb{R}^{n \times n}$, 由定理 17 可选正交阵 \mathbf{U} 使 $\mathbf{H} = \mathbf{U}^T \mathbf{A} \mathbf{U}$ 为上海森伯格阵, 对 \mathbf{H} 应用 QR 算法.

QR 算法: $\mathbf{H} = \mathbf{Q} \mathbf{R}$

对于 $k = 1, 2, \dots$

$$\begin{aligned}\mathbf{H}_k &= \mathbf{Q}_k \mathbf{R}_k && (\text{QR 分解}) \\ \mathbf{H}_{k+1} &= \mathbf{R}_k \mathbf{Q}_k\end{aligned}\quad (4.6)$$

不失一般性, 可假设由(4.6)迭代产生的每一个上海森伯格阵 \mathbf{H}_k 都是不可约的, 否则, 若在某步有

$$\mathbf{H}_{k+1} = \begin{array}{cc|c} \mathbf{H}_1 & \mathbf{H}_2 & p \\ \mathbf{0} & \mathbf{H}_2 & n-p \end{array}.$$

于是, 这个问题就分离为 \mathbf{H}_1 与 \mathbf{H}_2 两个较小的问题. 当 $p = n - 1$ 或 $n - 2$ 时, 有

$$\mathbf{H}_{k+1} = \begin{array}{cc|c} \mathbf{H}_1 & \mathbf{H}_2 & n-1 \\ \mathbf{0} & h_{n,n}^{(k+1)} & 1 \end{array}$$

或

$$\mathbf{H}_{k+1} = \begin{array}{cc|c} \mathbf{H}_1 & \mathbf{H}_2 & n-2 \\ \mathbf{0} & * & * \\ & * & 2 \end{array},$$

即可求出 \mathbf{H} 的特征值 $\lambda_n = h_{n,n}^{(k+1)}$ 或 λ_{n-1}, λ_n (由 \mathbf{H}_{k+1} 右下角二阶阵的特征值求得), 且求 \mathbf{H} 的其余特征值时, 转化为降阶求 \mathbf{H}_1 的特征值.

实际上, 每当 \mathbf{H}_{k+1} 的次对角元适当小时, 就可进行分离. 例如, 如果

$$|h_{p+1,p}| \ll (|h_{pp}| + |h_{p+1,p+1}|),$$

就把 $h_{p+1,p}$ 视为零. 一般取 $t = 10^{-t}$, 其中 t 是计算中有效数字的

位数 .

8.4.3 用单步 QR 方法计算上海森伯格阵特征值

上海森伯格阵的单步 QR 方法:选取 s_k 并设

$$\mathbf{H} = \begin{matrix} h_{11} & h_{12} & \dots & h_{1n} \\ h_{21} & h_{22} & \dots & h_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ h_{n,n-1} & h_{n,n} \end{matrix} = \mathbf{H} \quad (\text{设 } \mathbf{H} \text{ 为不可约阵}) .$$

对于 $k=1, 2, \dots$ (用位移来加速收敛)

$$\mathbf{H}_k - s_k \mathbf{I} = \mathbf{Q}_k \mathbf{R}_k$$

$$\mathbf{H}_{k+1} = \mathbf{R}_k \mathbf{Q}_k + s_k \mathbf{I}$$

由 $\mathbf{H}_k - \mathbf{H}_{k+1}$ 实际计算为

(1) 左变换: $\mathbf{P}_{n-1,n} \dots \mathbf{P}_{23} \mathbf{P}_{12} (\mathbf{H} - s_k \mathbf{I}) = \mathbf{R}$ (上三角阵) .

(2) 右变换: $\mathbf{H} = \mathbf{R} \mathbf{P}_{12}^T \mathbf{P}_{23}^T \dots \mathbf{P}_{n-1,n}^T + s_k \mathbf{I}$.

其中 $\mathbf{P}_{k,k+1} = \mathbf{P}(k, k+1)$ 为平面旋转阵 .

(1) 左变换计算

$$h_{kk} - s_k \quad (k = 1, 2, \dots, n),$$

确定平面旋转阵 $\mathbf{P}_{12} = \mathbf{P}(1, 2)$ 使

$$\begin{matrix} r_{11} & h_{12}^{(2)} & h_{13}^{(2)} & \dots & h_{1n}^{(2)} \\ 0 & h_{22}^{(2)} & h_{23}^{(2)} & \dots & h_{2n}^{(2)} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ h_{n,n-1} & h_{n,n} \end{matrix}$$

$$\mathbf{P}_{12} (\mathbf{H} - s_k \mathbf{I}) = \begin{matrix} 0 & h_{22} & h_{32} & \dots & h_{n2} \\ & \vdots & \vdots & \ddots & \vdots \\ & & & \dots & \end{matrix} .$$

设已完成第 1 次, … 第 $k-1$ 次左变换, 即有

$$\mathbf{P}_{k-1,k} \dots \mathbf{P}_{23} \mathbf{P}_{12} (\mathbf{H} - s_k \mathbf{I}) =$$

$$\begin{array}{cccccc}
 h_{11} & \dots & h_{1, k-1}^{(2)} & h_{1k}^{(2)} & \dots & h_{1n}^{(2)} \\
 W & \dots & \dots & \dots & \dots & \dots \\
 & & r_{k-1, k-1}^{(k)} & h_{k-1, k}^{(k)} & \dots & h_{k-1, n}^{(k)} \\
 & & h_{kk}^{(k)} & \dots & h_{kn}^{(k)} & \dots \\
 & & h_{k+1, k} & \dots & h_{k+1, n} & \\
 W & \dots & & & & \\
 & & h_{n, n-1} & h_{nn} & &
 \end{array} \quad (4.7)$$

确定平面旋转阵 $\mathbf{P}_{k, k+1} = \mathbf{P}(k, k+1)$, 使 $h_{k+1, k}$ 变为 0, 且完成第 k 次左变换 $\mathbf{P}_{k, k+1} (\mathbf{P}_{k-1, k} \dots \mathbf{P}_{12} (\mathbf{H} - s \mathbf{I}))$ 计算(只需计算(4.7)阵第 k 行及第 $k+1$ 行元素).

继续这一过程, 最后有

$$\mathbf{P}_{n-1, n} \dots \mathbf{P}_{12} (\mathbf{H} - s \mathbf{I}) = \mathbf{R} \quad (\text{上三角阵}).$$

(2) 右变换计算

$$\mathbf{H} = \mathbf{R} \mathbf{P}_{12}^T \mathbf{P}_{23}^T \dots \mathbf{P}_{n-1, n}^T + s \mathbf{I},$$

在第 k 次右变换 $(\mathbf{R} \mathbf{P}_{12}^T \dots) \mathbf{P}_{k, k+1}^T$ 中, 只需计算 $\mathbf{R} \mathbf{P}_{12}^T \dots \mathbf{P}_{k-1, k}^T$ 第 k 列及第 $k+1$ 列元素.

$$h_{k, k} = h_{k, k} + s \quad (k = 1, 2, \dots, n).$$

最后

$$\begin{aligned}
 \mathbf{H} &= \mathbf{R} \mathbf{P}_{12}^T \dots \mathbf{P}_{n-1, n}^T + s \mathbf{I} \\
 &\quad * * \dots * \\
 &\quad * * \dots * \\
 &= \begin{array}{ccc|c} & & & \\ W & W & \dots & \\ & * & * & \end{array} \quad (\text{为上海森伯格阵}).
 \end{aligned}$$

由上述讨论指出, 如果 $\mathbf{H} \in \mathbb{R}^{n \times n}$ 为上海森伯格阵, 则用 QR 算法产生的 $\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_k, \dots$ 亦是上海森伯格阵. 即上海森伯格阵在 QR 变换下形式不变.

下述定理讨论一个极端的情况

定理 22 设:(1) $\mathbf{H} \in \mathbb{R}^{n \times n}$ 为不可约上海森伯格阵; (2) μ 为

$\mathbf{H} = \mathbf{H}$ 一个特征值 则 QR 方法

$$\mathbf{H} - \mu \mathbf{I} = \mathbf{QR} \quad (\text{QR 分解})$$

$$\mathbf{H} = \mathbf{RQ} + \mu \mathbf{I}$$

中 $h_{n,n-1}^{(2)} = 0$, $h_{n,n}^{(2)} = \mu$.

证明 记

$$\mathbf{R} = \begin{matrix} r_{11} & \dots & r_{1n} \\ & \ddots & \dots \\ & & r_{nn} \end{matrix} \quad (\text{上三角阵}).$$

由设 \mathbf{H} 为不可约阵, 则上海森伯格阵 $\mathbf{H} - \mu \mathbf{I}$ 亦为不可约. 由将上海森伯格阵 $\mathbf{H} - \mu \mathbf{I}$ 约化为上三角阵 \mathbf{R} 的平面旋转变换的取法可知

$$|r_{ii}| / |h_{i+1,i}| / 0 \quad (i = 1, 2, \dots, n-1),$$

又因为 $\mathbf{Q}^T(\mathbf{H} - \mu \mathbf{I}) = \mathbf{R}$ 为奇异矩阵, 从而得到 $r_{nn} = 0$. 因此, \mathbf{H} 的最后一行为 $(0, 0, \dots, 0, \mu)$, 即

$$h_{n,n-1}^{(2)} = 0, \quad h_{n,n}^{(2)} = \mu.$$

这就启发我们在 QR 方法迭代中, 参数 s_k 可选为 $h_{n,n}^{(k)}$, 即 \mathbf{H}_k 的 (n, n) 元素. 通常可以作为特征值的最好近似.

算法 3(上海森伯格阵的 QR 算法) 给定 $\mathbf{H} \in \mathbb{R}^{n \times n}$ 为上海森伯格阵, 本算法计算

$$\mathbf{H} - s \mathbf{I} = \mathbf{Q} \mathbf{R} \quad (\text{QR 分解}) \quad (\text{取 } s = h_{n,n})$$

$$\mathbf{H} = \mathbf{R} \mathbf{Q} + s \mathbf{I}$$

且 \mathbf{H} 覆盖 $\mathbf{H} (\mathbf{H} = \mathbf{H})$

$$1. \quad h_{11} - s$$

2. 对于 $k = 1, 2, \dots, n-1$

$$(1) \quad h_{k+1,k+1} - h_{k+1,k+1} - s$$

(2) 确定 $\mathbf{P}(k, k+1)$ 使

$$\begin{matrix} \alpha & s_k & h_{kk} & = & r_{kk} \\ -s_k & c_k & h_{k+1,k} & = & 0 \end{matrix}$$

(3) 左变换

对于 $j = k, \dots, n$

$$\begin{array}{cccc} h_{kj} & & c_k & s_k \\ & & -s_k & c_k \\ h_{k+1,j} & & & h_{k+1,j} \end{array}$$

3. 对于 $k = 1, 2, \dots, n-1$

(1) 右变换

对于 $i = 1, 2, \dots, k+1$

$$\begin{array}{ccc} (h_{ik}, h_{i,k+1}) & & c_k \\ & & -s_k \\ & & s_k \\ & & c_k \end{array}$$

$$(2) h_{kk} \quad h_{kk} + s$$

$$4. h_{n,n} \quad h_{n,n} + s$$

如果用不同的位移 $s_k = h_{n,n}^{(k)}$, 反复应用算法 3 就产生正交相似的上海森伯格阵序列 $\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_k, \dots$. 当 $h_{n,n}^{(k)}$ 充分小时, 可将它置为零就得到 \mathbf{H} 的近似特征值 $\lambda_n = h_{n,n}^{(k)}$. 再将矩阵降阶, 对较小矩阵连续应用算法 .

例 8 用 QR 方法计算对称三对角矩阵

$$\mathbf{A} = \begin{pmatrix} 2 & 1 & 0 \\ 1 & 3 & 1 \\ 0 & 1 & 4 \end{pmatrix}.$$

的全部特征值 .

解 选取 $s_1 = d_{11}^{(1)}$, 则 $s_1 = 4$.

$$\mathbf{P}_{23} \mathbf{P}_{12} (\mathbf{A} - s_1 \mathbf{I}) = \mathbf{R}$$

$$\begin{aligned} & 2.2361 \quad -1.342 \quad 0.4472 \\ & = \quad \quad \quad 1.0954 \quad -0.3651 \quad , \\ & \quad \quad \quad \quad \quad \quad 0.81650 \end{aligned}$$

$$\mathbf{A}_1 = \mathbf{R} \mathbf{P}_{12}^T \mathbf{P}_{23}^T + s_1 \mathbf{I}$$

$$\begin{aligned}
 & 1 .4000 \quad 0 .4899 \quad 0 \\
 = & 0 .4899 \quad 3 .2667 \quad 0 .7454 \\
 & 0 \quad 0 .7454 \quad 4 .3333 \\
 & 1 .2915 \quad 0 .2017 \quad 0 \\
 \mathbf{A}_3 = & 0 .2017 \quad 3 .0202 \quad 0 .2724 \\
 & 0 \quad 0 .2724 \quad 4 .6884 \\
 & 1 .2737 \quad 0 .0993 \quad 0 \\
 \mathbf{A}_4 = & 0 .0993 \quad 2 .9943 \quad 0 .0072 \\
 & 0 \quad 0 .0072 \quad 4 .7320 \\
 & 1 .2694 \quad 0 .0498 \quad 0 \\
 \mathbf{A}_5 = & 0 .0498 \quad 2 .9986 \quad 0 \\
 & 0 \quad 0 \quad 4 .7321 \\
 \mathbf{H} = & 1 .2694 \quad 0 .0498 \\
 & 0 .0498 \quad 2 .9986
 \end{aligned}$$

现在收缩,继续对 \mathbf{A}_5 的子矩阵 \mathbf{H} $\in \mathbb{R}^{2 \times 2}$ 进行变换,得到

$$\mathbf{H} = \mathbf{P}_{12} (\mathbf{H} - s_k \mathbf{I}) \mathbf{P}_{12}^T + s_k \mathbf{I} = \begin{matrix} 1 .2680 & - 4 \times 10^{-5} \\ - 4 \times 10^{-5} & 3 .0000 \end{matrix},$$

故求得 \mathbf{A} 近似特征值为

$$3 .7321, \quad 2 .0000, \quad 1 .2680.$$

而 \mathbf{A} 的特征值是

$$3 = 3 + 3 .7321, \quad 2 = 3 .0, \quad 1 = 3 - 3 .2679.$$

算法 3 是在实数中进行选择位移 $s_k = h_{n,n}^{(k)}$, 不能逼近一个复特征值,所以算法 3 不能用来计算 \mathbf{H} 的复特征值.

8.4.4* 双步 QR 方法(隐式 QR 方法)

在本章第 3 节中将 $\mathbf{A} \in \mathbb{R}^{n \times n}$ 经过正交相似变换化为上海森伯格矩阵 \mathbf{H} , 即 $\mathbf{U}^T \mathbf{A} \mathbf{U} = \mathbf{H}$, 其中 \mathbf{H} 不是唯一的.但是,如果规定了正交矩阵 \mathbf{U} 的第一列,则 \mathbf{U} 和 \mathbf{H} 除差 ± 1 因子外唯一.

定理 23(隐式 Q 定理) 设 $\mathbf{A} \in \mathbb{R}^{n \times n}$, 且:

- (1) $\mathbf{Q} = (\mathbf{q} \ \mathbf{q} \dots \mathbf{q})$ 及 $\mathbf{V} = (\mathbf{v} \ \mathbf{v} \dots \mathbf{v})$ 都是正交阵, 且有 $\mathbf{Q}^T \mathbf{A} \mathbf{Q} = \mathbf{H}$, $\mathbf{V}^T \mathbf{A} \mathbf{V} = \mathbf{G}$ 都是上海森伯格阵 .
- (2) \mathbf{H} 为不可约上海森伯格阵, 且 $\mathbf{q} = \mathbf{v}$ (即 \mathbf{Q} 与 \mathbf{V} 第 1 列相同) . 则:

- (1) $\mathbf{v} = \pm \mathbf{q}$, 且 $|h_{i, i+1}| = |g_{i, i+1}|$ ($i = 2, \dots, n$);
- (2) $\mathbf{G} = \mathbf{D}^{-1} \mathbf{H} \mathbf{D}$, 其中 $\mathbf{D} = \text{diag}(1, \pm 1, \dots, \pm 1)$, 即 \mathbf{H} 和 \mathbf{G} 在 $\mathbf{G} = \mathbf{D}^{-1} \mathbf{H} \mathbf{D}$ 意义上“本质上相等”.

算法 3 不能用来求 \mathbf{H} 的一个复特征值, 当 \mathbf{H} (上海森伯格阵) 的依模最小特征值是复数时, 位移参数 s_k, s_{k+1} 可取为某步 \mathbf{H}_k 右下角的二阶矩阵

$$\mathbf{G} = \begin{matrix} h_{n-1, n-1} & h_{n-1, n} \\ h_{n, n-1} & h_{n, n} \end{matrix} \quad (4.8)$$

的特征值 .

当 \mathbf{G} 的特征值 s_1 与 s_2 为复数时, 如果应用算法 3 就要引进复数运算, 这对于实矩阵 \mathbf{H} 是不必要的, 事实上, 在某些条件下, 可以用正交相似变换将 \mathbf{H} 约化为实 Schur 型 .

下面引进隐式位移的 QR 方法, 即用 s_1 与 s_2 作位移连续进行二次单步的 QR 迭代, 使用复位移, 又避免复数运算 .

(1) 设 $\mathbf{H} = \mathbf{H} \in \mathbb{R}^{n \times n}$ 为上海森伯格阵, 取共轭复数 s_1, s_2 作两步位移的 QR 方法, 即

$$\begin{aligned} \mathbf{H} - s_1 \mathbf{I} &= \mathbf{Q} \mathbf{R}, \\ \mathbf{H} &= \mathbf{R} \mathbf{Q} + s_1 \mathbf{I} = \mathbf{Q}^T \mathbf{H} \mathbf{Q}, \\ \mathbf{H} - s_2 \mathbf{I} &= \mathbf{Q} \mathbf{R}, \\ \mathbf{H} &= \mathbf{R} \mathbf{Q} + s_2 \mathbf{I} = \mathbf{Q}^T \mathbf{Q}^T \mathbf{H} \mathbf{Q} \mathbf{Q} = \mathbf{Q}^T \mathbf{H} \mathbf{Q}, \end{aligned} \quad (4.9)$$

其中, $\mathbf{Q} = \mathbf{Q} \mathbf{Q}$, $\mathbf{R} = \mathbf{R} \mathbf{R}$.

显然 $\mathbf{M} = (\mathbf{H} - s_1 \mathbf{I})(\mathbf{H} - s_2 \mathbf{I})$ 有 QR 分解

$$\mathbf{M} = \mathbf{Q} \mathbf{R}. \quad (4.10)$$

事实上,由(4.9)式并利用 $\mathbf{H} - s\mathbf{I} = \mathbf{Q}^T (\mathbf{H} - s\mathbf{I}) \mathbf{Q} = \mathbf{Q} \mathbf{R}$ 有

$$\begin{aligned}\mathbf{M} &= (\mathbf{H} - s\mathbf{I}) \mathbf{Q} \mathbf{R} \\ &= (\mathbf{Q} \mathbf{Q} \mathbf{R} \mathbf{Q}^T) \mathbf{Q} \mathbf{R} \\ &= \mathbf{Q} \mathbf{Q} \mathbf{R} \mathbf{R} = \mathbf{QR}.\end{aligned}$$

且 \mathbf{M} 阵为实矩阵,这是因为(即使 \mathbf{G} 特征值为复数)

$$\mathbf{M} = \mathbf{H}^2 - (s_1 + s_2) \mathbf{H} + s_1 s_2 \mathbf{I}, \quad (4.11)$$

其中 $s_1 + s_2 = h_{n-1,n-1} + h_{nn} = s$, $s_1 s_2 = h_{n-1,n-1} h_{nn} - h_{n,n-1} h_{n-1,n} = t$ 为实数.于是,(4.10)式为实矩阵 \mathbf{M} 的 QR 分解,并且可以选取 \mathbf{Q} 和 \mathbf{Q} 使 $\mathbf{Q} = \mathbf{Q} \mathbf{Q}$ 为实的正交阵.由此得出

$$\mathbf{H} = (\mathbf{Q} \mathbf{Q})^T \mathbf{H} (\mathbf{Q} \mathbf{Q}) = \mathbf{Q}^T \mathbf{H} \mathbf{Q}$$

是实矩阵.

如果用下述算法就能保证 \mathbf{H} 是实矩阵

- (a) 直接形成实矩阵 $\mathbf{M} = \mathbf{H}^2 - s\mathbf{H} + t\mathbf{I}$
- (b) 计算 \mathbf{M} 阵的实 QR 分解 $\mathbf{M} = \mathbf{QR}$
- (c) 令 $\mathbf{H}_3 = \mathbf{Q}^T \mathbf{H} \mathbf{Q}$

但是(a)需要 $O(n^3)$ 次乘法运算,不实用.

(2) 根据隐式 Q 定理,如果按下述算法进行,就有可能用 $O(n^2)$ 次运算来实现从 \mathbf{H} 到 \mathbf{H}_3 的转换.

- (a) 求与 \mathbf{Q} 有相同第一列的正交阵 \mathbf{P}_0
- (b) 应用豪斯霍尔德方法将 $\mathbf{P}_0^T \mathbf{H} \mathbf{P}_0$ 化为一个上三角形阵,即

$$\mathbf{P}_{n-2} \dots \mathbf{P}_2 \mathbf{P}_1 (\mathbf{P}_0^T \mathbf{H} \mathbf{P}_0) \mathbf{P}_1 \mathbf{P}_2 \dots \mathbf{P}_{n-2} = \mathbf{H}.$$

记 $\mathbf{Q} = \mathbf{P}_0 \mathbf{P}_1 \dots \mathbf{P}_{n-2}$, 上式为

$$\mathbf{Q}^T \mathbf{H} \mathbf{Q} = \mathbf{H}.$$

显然, \mathbf{Q} 的第一列与 \mathbf{P}_0 的第一列相同,即 \mathbf{Q} 与 \mathbf{Q} 第一列相同 ($\mathbf{Q} \mathbf{e} = \mathbf{P}_0 \mathbf{e} = \mathbf{Q} \mathbf{e}$).若 $\mathbf{Q}^T \mathbf{H} \mathbf{Q}$ 与 $\mathbf{Q}^T \mathbf{H} \mathbf{Q}$ 两者都是不可约上三角形阵,则由隐式 Q 定理 \mathbf{H} 与 \mathbf{H}_3 本质上相等.

- (3) 如何寻求正交阵 \mathbf{P}_0 .

由于 $\mathbf{M} = \mathbf{QR}$ (为 \mathbf{M} 的 QR 分解), 则

$$\mathbf{Me} = \mathbf{QRe} = n_1 \mathbf{Qe}.$$

这说明 \mathbf{Q} 的第一列即是 \mathbf{M} 第一列的一个倍数, 于是, 对 \mathbf{M} 阵的第一列(非零)寻求初等反射阵 \mathbf{P}_0 使

$$\mathbf{P}_0(\mathbf{Me}) = n_1 \mathbf{e} \quad (\text{其中 } n_1 = -)$$

即

$$\mathbf{Me} = n_1 \mathbf{P}_0 \mathbf{e}.$$

这说明 \mathbf{P}_0 与 \mathbf{Q} 具有相同的第一列.

由于 $\mathbf{M} = (\mathbf{H} - s_1 \mathbf{I})(\mathbf{H} - s_2 \mathbf{I})$, 则

$$\mathbf{Me} = (x, y, z, 0, \dots, 0)^T$$

其中

$$x = (h_{11} - s_1)(h_{11} - s_2) + h_{12} h_{21}$$

$$= h_{11}^2 + h_{12} h_{21} - s_1 h_{11} + t,$$

$$y = (h_{11} - s_2)h_{21} + (h_{22} - s_1)h_{11} \quad (4.12)$$

$$= h_{11}(h_{11} + h_{22} - s_1),$$

$$z = h_{21} h_{22}.$$

双步 QR 方法: 设 $\mathbf{H} = \mathbf{H} \in \mathbf{R}^{n \times n}$ 为不可约上三角阵.

(a) 计算 \mathbf{M} 阵的第一列. 即按(4.12)式计算

$$\mathbf{Me} = (x, y, z, 0, \dots, 0)^T;$$

(b) 确定初等反射阵 \mathbf{P}_0 使

$$\mathbf{P}_0(\mathbf{Me}) = -\mathbf{e},$$

x

即确定初等反射阵 $\mathbf{R} \in \mathbf{R}^{3 \times 3}$ 使 $\mathbf{R} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = -\mathbf{e}$;

$$\mathbf{P}_0 = \begin{matrix} \mathbf{R} & & \\ & & 3 \\ & & \mathbf{I}_{n-3} \end{matrix};$$

(c) 计算初等反射阵 $\mathbf{P}_1, \mathbf{P}_2, \dots, \mathbf{P}_{n-2}$ 使

$$\mathbf{P}_{n-2} \dots \mathbf{P}_2 \mathbf{P}_1 (\mathbf{P}_0 \mathbf{H} \mathbf{R}) \mathbf{P}_1 \mathbf{P}_2 \dots \mathbf{P}_{n-2} = \mathbf{H}$$

为上三角阵, 则 $\mathbf{Q} = \mathbf{Q}_1 \mathbf{Q}_2 \dots \mathbf{Q}_{n-2}$ 与 $\mathbf{Q} = \mathbf{P}_0 \mathbf{P}_1 \dots \mathbf{P}_{n-2}$ 第一列相同且 $\mathbf{H} = \mathbf{H}_0$.

这样上面的算法就完成了从 \mathbf{H} 到 \mathbf{H}' 的变换, 但没有明显地应用到位移 s_1 和 s_2 .

算法的具体实现可参看文献[1]及[19], 本算法在数学库中均有相应软件, 可直接使用.

评 注

本章介绍了计算一般矩阵主特征值及对应特征向量的迭代法(幂法), 这种方法在计算过程中原始矩阵 \mathbf{A} 始终不变. 因此, 这种方法适于求高阶稀疏矩阵的特征值问题. 主特征值为复特征值的情况可参考文献[3].

反幂法是计算矩阵特征向量的一个有效方法, 主要用来计算海森伯格阵或三对角矩阵的对应于一个给定的近似特征值的特征向量.

本章还介绍了用正交相似变换的方法约化一般矩阵 \mathbf{A} 为上三角形的豪斯霍尔德方法, 计算一般矩阵 \mathbf{A} 全部特征值问题的 QR 方法. 这些方法都是利用矩阵的正交相似变换约化矩阵 \mathbf{A} 为某种简单形式, 以期求 \mathbf{A} 的特征值的方法. 这种方法将破坏原始矩阵. QR 方法(与豪斯霍尔德方法结合使用)具有收敛快, 精度高的特点, 是求任意实矩阵(中, 小型)特征值问题的最有效的方法之一. 对计算矩阵特征值问题有兴趣的读者可进一步参考文献[13], [19], [20].

习 题

1. 用幂法计算下列矩阵的主特征值及对应的特征向量:

$$(a) \mathbf{A}_1 = \begin{matrix} 7 & 3 & -2 \\ 3 & 4 & -1 \\ -2 & -1 & 3 \end{matrix}, \quad (b) \mathbf{A}_2 = \begin{matrix} 3 & -4 & 3 \\ -4 & 6 & 3 \\ 3 & 3 & 1 \end{matrix},$$

当特征值有 3 位小数稳定时迭代终止 .

2 . 利用反幂法求矩阵

$$\begin{matrix} 6 & 2 & 1 \\ 2 & 3 & 1 \\ 1 & 1 & 1 \end{matrix}$$

的最接近于 6 的特征值及对应的特征向量 .

3 . 求矩阵

$$\begin{matrix} 4 & 0 & 0 \\ 0 & 3 & 1 \\ 0 & 1 & 3 \end{matrix}$$

与特征值 4 对应的特征向量 .

4 . a) 设 \mathbf{A} 是对称矩阵, 和 \mathbf{x} ($\|\mathbf{x}\|_2 = 1$) 是 \mathbf{A} 的一个特征值及相应的特征向量 . 又设 \mathbf{P} 为一个正交阵, 使

$$\mathbf{P}\mathbf{x} = \mathbf{e} = (1, 0, \dots, 0)^T .$$

证明 $\mathbf{B} = \mathbf{P}\mathbf{A}\mathbf{P}^T$ 的第一行和第一列除了 外其余元素均为零 .

(b) 对于矩阵

$$\mathbf{A} = \begin{matrix} 2 & 10 & 2 \\ 10 & 5 & -8 \\ 2 & -8 & 11 \end{matrix},$$

$= 9$ 是其特征值, $\mathbf{x} = \frac{2}{3}, \frac{1}{3}, \frac{2}{3}^T$ 是相应于 9 的特征向量, 试求一初等反射阵 \mathbf{P} , 使 $\mathbf{P}\mathbf{x} = \mathbf{e}$, 并计算 $\mathbf{B} = \mathbf{P}\mathbf{A}\mathbf{P}^T$.

5 . 利用初等反射阵将

$$\mathbf{A} = \begin{matrix} 1 & 3 & 4 \\ 3 & 1 & 2 \\ 4 & 2 & 1 \end{matrix}$$

正交相似约化为对称三对角阵 .

6 . 设 \mathbf{A}_{n-1} 是由豪斯霍尔德方法得到的矩阵, 又设 \mathbf{y} 是 \mathbf{A}_{n-1} 的一个特征向量 .

(a) 证明矩阵 \mathbf{A} 对应的特征向量是 $\mathbf{x} = \mathbf{P}_1 \mathbf{P}_2 \dots \mathbf{P}_{n-2} \mathbf{y}$,

(b) 对于给出的 \mathbf{y} 应如何计算 \mathbf{x} ?

7. 用带位移的 QR 方法计算

$$(a) \mathbf{A} = \begin{pmatrix} 1 & 2 & 0 \\ 2 & -1 & 1 \\ 0 & 1 & 3 \end{pmatrix}, \quad (b) \mathbf{B} = \begin{pmatrix} 3 & 1 & 0 \\ 1 & 2 & 1 \\ 0 & 1 & 1 \end{pmatrix}$$

的全部特征值 .

8. 试用初等反射阵将

$$\mathbf{A} = \begin{pmatrix} 1 & 1 & 1 \\ 2 & -1 & -1 \\ 2 & -4 & 5 \end{pmatrix}$$

分解为 QR 的形式 , 其中 \mathbf{Q} 为正交阵 , \mathbf{R} 为上三角阵 .

9. 设

$$\mathbf{A} = \begin{pmatrix} 4 & -1 & & \\ -1 & 4 & -1 & \\ & & \mathbf{W} & \mathbf{W} & \mathbf{W} \\ & & -1 & 4 & -1 \\ & & -1 & 4 & \end{pmatrix} \in \mathbf{R}^{n \times n},$$

试确定 \mathbf{A} 及 \mathbf{A}^{-1} 特征值的界 .

10. 设 $\mathbf{A} = \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} & 3 \\ \mathbf{0} & \mathbf{A}_{22} & 2 \end{pmatrix}$, 又设 i 为 \mathbf{A}_{11} 的特征值 , j 为 \mathbf{A}_{22} 的特征值 ,

$\mathbf{x} = (1, 2, 3)^T$ 为对应于 i , \mathbf{A}_{11} 的特征向量 , $\mathbf{y} = (1, 2)^T$ 为对应于 j , \mathbf{A}_{22} 的特征向量 . 求证 :

(1) i, j 为 \mathbf{A} 的特征值 .

(2) $\mathbf{x} = (1, 2, 3, 0, 0)^T$ 为对应于 i , \mathbf{A} 的特征向量 ,

$\mathbf{y} = (0, 0, 0, 1, 2)^T$ 为对应于 j , \mathbf{A} 的特征向量 .

第9章 常微分方程初值问题数值解法

9.1 引言

科学技术中常常需要求解常微分方程的定解问题.这类问题最简单的形式,是本章将要着重考察的一阶方程的初值问题

$$y = f(x, y), \quad (1.1)$$

$$y(x_0) = y_0. \quad (1.2)$$

我们知道,只要函数 $f(x, y)$ 适当光滑——譬如关于 y 满足利普希茨(Lipschitz)条件

$$|f(x, y) - f(x, \bar{y})| \leq L |y - \bar{y}|. \quad (1.3)$$

理论上就可以保证初值问题(1.1),(1.2)的解 $y = y(x)$ 存在并且唯一.

虽然求解常微分方程有各种各样的解析方法,但解析方法只能用来求解一些特殊类型的方程,实际问题中归结出来的微分方程主要靠数值解法.

所谓数值解法,就是寻求解 $y(x)$ 在一系列离散节点

$$x_1 < x_2 < \dots < x_n < x_{n+1} < \dots$$

上的近似值 $y_1, y_2, \dots, y_n, y_{n+1}, \dots$.相邻两个节点的距离 $h_i = x_{i+1} - x_i$ 称为步长.今后如不特别说明,总是假定 $h_i = h$ ($i = 1, 2, \dots$) 为定数,这时节点为 $x_n = x_0 + nh, n = 0, 1, 2, \dots$.

初值问题(1.1),(1.2)的数值解法有个基本特点,它们都采取“步进式”,即求解过程顺着节点排列的次序一步一步地向前推进.描述这类算法,只要给出用已知信息 $y_n, y_{n-1}, y_{n-2}, \dots$ 计算 y_{n+1} 的递推公式.

首先,要对方程(1.1)离散化,建立求数值解的递推公式.一类是计算 y_{n+1} 时只用到前一点的值 y_n ,称为单步法.另一类是用到 y_{n+1} 前面 k 点的值 $y_n, y_{n-1}, \dots, y_{n-k+1}$,称为 **k 步法**.其次,要研究公式的局部截断误差和阶,数值解 y_n 与精确解 $y(x_n)$ 的误差估计及收敛性,还有递推公式的计算稳定性等问题.

9.2 简单的数值方法与基本概念

9.2.1 欧拉法与后退欧拉法

我们知道,在 xy 平面上,微分方程(1.1)的解 $y = y(x)$ 称作它的积分曲线.积分曲线上一点 (x, y) 的切线斜率等于函数 $f(x, y)$ 的值.如果按函数 $f(x, y)$ 在 xy 平面上建立一个方向场,那么,积分曲线上每一点的切线方向均与方向场在该点的方向相一致.

基于上述几何解释,我们从初始点 $P_0(x_0, y_0)$ 出发,先依方向场在该点的方向推进到 $x = x_1$ 上一点 P_1 ,然后再从 P_1 依方向场的方向推进到 $x = x_2$ 上一点 P_2 ,循此前进做出一条折线 $\overline{P_0 P_1 P_2 \dots}$ (图9-1).

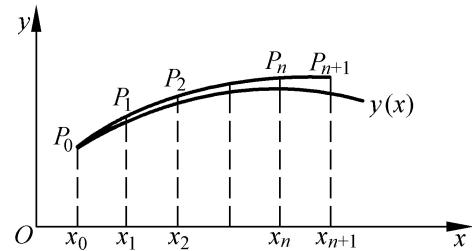


图 9-1

一般地,设已做出该折线的顶点 P_n ,过 $P_n(x_n, y_n)$ 依方向场的方向再推进到 $P_{n+1}(x_{n+1}, y_{n+1})$,显然两个顶点 P_n, P_{n+1} 的坐标有关系

$$\frac{y_{n+1} - y_n}{x_{n+1} - x_n} = f(x_n, y_n),$$

即

$$y_{n+1} = y_n + h f(x_n, y_n). \quad (2.1)$$

这就是著名的欧拉(Euler)公式.若初值 y_0 已知,则依公式(2.1)可逐步算出

$$y_1 = y_0 + hf(x_0, y_0),$$

$$y_2 = y_1 + hf(x_1, y_1),$$

...

例 1 求解初值问题

$$y = y - \frac{2x}{y} \quad (0 < x < 1), \quad (2.2)$$

$$y(0) = 1.$$

解 为便于进行比较, 本章将用多种数值方法求解上述初值问题. 这里先用欧拉方法, 欧拉公式的具体形式为

$$y_{n+1} = y_n + h \cdot y_n - \frac{2x_n}{y_n}.$$

取步长 $y=0.1$, 计算结果见表 9-1.

表 9-1 计算结果对比

x_n	y_n	$y(x_n)$	x_n	y_n	$y(x_n)$
0.1	1.1000	1.0954	0.6	1.5090	1.4832
0.2	1.1918	1.1832	0.7	1.5803	1.5492
0.3	1.2774	1.2649	0.8	1.6498	1.6125
0.4	1.3582	1.3416	0.9	1.7178	1.6733
0.5	1.4351	1.4142	1.0	1.7848	1.7321

初值问题(2.2)有解 $y = 1 + 2x$, 按这个解析式子算出的准确值 $y(x_n)$ 同近似值 y_n 一起列在表 9-1 中, 两者相比较可以看出欧拉方法的精度很差.

还可以通过几何直观来考察欧拉方法的精度. 假设 $y_n = y(x_n)$, 即顶点 P_n 落在积分曲线 $y = y(x)$ 上, 那么, 按欧拉方法做出的折线 $P_n P_{n+1}$ 便是 $y = y(x)$ 过点 P_n 的切线(图 9-2). 从图形上看, 这样定

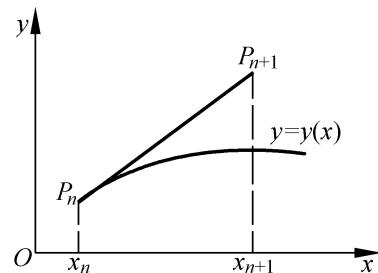


图 9-2

出的顶点 P_{n+1} 显著地偏离了原来的积分曲线, 可见欧拉方法是相当粗糙的.

为了分析计算公式的精度, 通常可用泰勒展开将 $y(x_{n+1})$ 在 x_n 处展开, 则有

$$\begin{aligned} y(x_{n+1}) &= y(x_n + h) \\ &= y(x_n) + y'(x_n)h + \frac{h^2}{2}y''(x_n), \quad n \in (x_n, x_{n+1}). \end{aligned}$$

在 $y_n = y(x_n)$ 的前提下, $f(x_n, y_n) = f(x_n, y(x_n)) = y'(x_n)$. 于是可得欧拉法(2.1)的公式误差

$$y(x_{n+1}) - y_{n+1} = \frac{h^2}{2}y''(x_n) - \frac{h^2}{2}y''(x_n), \quad (2.3)$$

称为此方法的局部截断误差.

如果对方程(1.1)从 x_n 到 x_{n+1} 积分, 得

$$y(x_{n+1}) = y(x_n) + \int_{x_n}^{x_{n+1}} f(t, y(t)) dt. \quad (2.4)$$

右端积分用左矩形公式 $hf(x_n, y(x_n))$ 近似, 再以 y_n 代替 $y(x_n)$, y_{n+1} 代替 $y(x_{n+1})$ 也得到(2.1), 局部截断误差也是(2.3).

如果在(2.4)中右端积分用右矩形公式 $hf(x_{n+1}, y(x_{n+1}))$ 近似, 则得另一个公式

$$y_{n+1} = y_n + hf(x_{n+1}, y_{n+1}), \quad (2.5)$$

称为后退的欧拉法.

后退的欧拉公式与欧拉公式有着本质的区别, 后者是关于 y_{n+1} 的一个直接的计算公式, 这类公式称作是显式的; 然而公式(2.5)的右端含有未知的 y_{n+1} , 它实际上是关于 y_{n+1} 的一个函数方程, 这类公式称作是隐式的.

显式与隐式两类方法各有特点. 考虑到数值稳定性等其他因素, 人们有时需要选用隐式方法, 但使用显式算法远比隐式方便.

隐式方程(2.5)通常用迭代法求解, 而迭代过程的实质是逐步显示化.

设用欧拉公式

$$y_{n+1}^{(0)} = y_n + hf(x_n, y_n)$$

给出迭代初值 $y_{n+1}^{(0)}$, 用它代入(2.5)式的右端, 使之转化为显式, 直接计算得

$$y_{n+1}^{(1)} = y_n + hf(x_{n+1}, y_{n+1}^{(0)}),$$

然后再用 $y_{n+1}^{(1)}$ 代入(2.5)式, 又有

$$y_{n+1}^{(2)} = y_n + hf(x_{n+1}, y_{n+1}^{(1)}).$$

如此反复进行, 得

$$y_{n+1}^{(k+1)} = y_n + hf(x_{n+1}, y_{n+1}^{(k)}) \quad (k = 0, 1, \dots). \quad (2.6)$$

由于 $f(x, y)$ 对 y 满足利普希茨条件(1.3). 由(2.6)减(2.5)得

$$\begin{aligned} |y_{n+1}^{(k+1)} - y_{n+1}| &= h |f(x_{n+1}, y_{n+1}^{(k)}) - f(x_{n+1}, y_{n+1})| \\ &\leq hL |y_{n+1}^{(k)} - y_{n+1}|. \end{aligned}$$

由此可知, 只要 $hL < 1$ 迭代法(2.6)就收敛到解 y_{n+1} . 关于后退欧拉方法的公式误差, 从积分公式看到它与欧拉法是相似的.

9.2.2 梯形方法

为得到比欧拉法精度高的计算公式, 在等式(2.4)右端积分中若用梯形求积公式近似, 并用 y_n 代替 $y(x_n)$, y_{n+1} 代替 $y(x_{n+1})$, 则得

$$y_{n+1} = y_n + \frac{h}{2} [f(x_n, y_n) + f(x_{n+1}, y_{n+1})], \quad (2.7)$$

称为梯形方法.

梯形方法是隐式单步法, 可用迭代法求解. 同后退的欧拉方法一样, 仍用欧拉方法提供迭代初值, 则梯形法的迭代公式为

$$\begin{aligned} y_{n+1}^{(0)} &= y_n + hf(x_n, y_n); \\ y_{n+1}^{(k+1)} &= y_n + \frac{h}{2} [f(x_n, y_n) + f(x_{n+1}, y_{n+1}^{(k)})] \quad (2.8) \\ (k &= 0, 1, 2, \dots). \end{aligned}$$

为了分析迭代过程的收敛性, 将(2.7)式与(2.8)式相减, 得

$$y_{n+1} - y_{n+1}^{(k+1)} = \frac{h}{2} [f(x_{n+1}, y_{n+1}) - f(x_{n+1}, y_{n+1}^{(k)})],$$

于是有

$$|y_{n+1} - y_{n+1}^{(k+1)}| \leq \frac{hL}{2} |y_{n+1} - y_{n+1}^{(k)}|,$$

式中 L 为 $f(x, y)$ 关于 y 的利普希茨常数. 如果选取 h 充分小, 使得

$$\frac{hL}{2} < 1,$$

则当 k 时有 $y_{n+1}^{(k)} = y_{n+1}$, 这说明迭代过程(2.8)是收敛的.

9.2.3 单步法的局部截断误差与阶

初值问题(1.1), (1.2)的单步法可用一般形式表示为

$$y_{n+1} = y_n + h (x_n, y_n, y_{n+1}, h), \quad (2.9)$$

其中多元函数 与 $f(x, y)$ 有关, 当 含有 y_{n+1} 时, 方法是隐式的, 若不含 y_{n+1} 则为显式方法, 所以显式单步法可表示为

$$y_{n+1} = y_n + h (x_n, y_n, h), \quad (2.10)$$

(x, y, h) 称为增量函数, 例如对欧拉法(2.1)有

$$(x, y, h) = f(x, y).$$

它的局部截断误差已由(2.3)给出, 对一般显式单步法则可如下定义.

定义 1 设 $y(x)$ 是初值问题(1.1), (1.2)的准确解, 称

$$T_{n+1} = y(x_{n+1}) - y(x_n) - h (x_n, y(x_n), h) \quad (2.11)$$

为显式单步法(2.10)的局部截断误差.

T_{n+1} 之所以称为局部的, 是假设在 x_n 前各步没有误差. 当 $y_n = y(x_n)$ 时, 计算一步, 则有

$$\begin{aligned} y(x_{n+1}) - y_{n+1} &= y(x_{n+1}) - [y_n + h (x_n, y_n, h)] \\ &= y(x_{n+1}) - y(x_n) - h (x_n, y(x_n), h) = T_{n+1}. \end{aligned}$$

所以,局部截断误差可理解为用方法(2.10)计算一步的误差,也即公式(2.10)中用准确解 $y(x)$ 代替数值解产生的公式误差.根据定义,显然欧拉法的局部截断误差

$$\begin{aligned} T_{n+1} &= y(x_{n+1}) - y(x_n) - hf(x_n, y(x_n)) \\ &= y(x_n + h) - y(x_n) - hy(x_n) \\ &= \frac{h^2}{2} y''(x_n) + O(h^3), \end{aligned}$$

即为(2.3)的结果.这里 $\frac{h^2}{2} y''(x_n)$ 称为局部截断误差主项.显然

$T_{n+1} = O(h^2)$.一般情形的定义如下.

定义2 设 $y(x)$ 是初值问题(1.1),(1.2)的准确解,若存在最大整数 p 使显式单步法(2.10)的局部截断误差满足

$$T_{n+1} = y(x + h) - y(x) - h(x, y, h) = O(h^{p+1}), \quad (2.12)$$

则称方法(2.10)具有 p 阶精度.

若将(2.12)展开式写成

$$T_{n+1} = (x_n, y(x_n)) h^{p+1} + O(h^{p+2}),$$

则 $(x_n, y(x_n)) h^{p+1}$ 称为局部截断误差主项.

以上定义对隐式单步法(2.9)也是适用的.例如,对后退欧拉法(2.5)其局部截断误差为

$$\begin{aligned} T_{n+1} &= y(x_{n+1}) - y(x_n) - hf(x_{n+1}, y(x_{n+1})) \\ &= hy(x_n) + \frac{h^2}{2} y''(x_n) + O(h^3) \\ &\quad - h[y(x_n) + hy(x_n) + O(h^2)] \\ &= -\frac{h^2}{2} y''(x_n) + O(h^3). \end{aligned}$$

这里 $p=1$,是1阶方法,局部截断误差主项为 $-\frac{h^2}{2} y''(x_n)$.

同样对梯形法(2.7)有

$$\begin{aligned}
 T_{n+1} &= y(x_{n+1}) - y(x_n) - \frac{h}{2}[y(x_n) + y(x_{n+1})] \\
 &= hy(x_n) + \frac{h^2}{2}y'(x_n) + \frac{h^3}{3!}y''(x_n) \\
 &\quad - \frac{h}{2}[y(x_n) + y(x_n) + hy(x_n) + \frac{h^2}{2}y'(x_n)] + O(h^4) \\
 &= -\frac{h^3}{12}y''(x_n) + O(h^4).
 \end{aligned}$$

所以梯形方法(2.7)是二阶的, 其局部误差主项为 $-\frac{h^3}{12}y''(x_n)$.

9.2.4 改进的欧拉公式

我们看到, 梯形方法虽然提高了精度, 但其算法复杂, 在应用迭代公式(2.9)进行实际计算时, 每迭代一次, 都要重新计算函数 $f(x, y)$ 的值, 而迭代又要反复进行若干次, 计算量很大, 而且往往难以预测. 为了控制计算量, 通常只迭代一两次就转入下一步的计算, 这就简化了算法.

具体地说, 我们先用欧拉公式求得一个初步的近似值 \tilde{y}_{n+1} , 称之为预测值, 预测值 \tilde{y}_{n+1} 的精度可能很差, 再用梯形公式(2.7)将它校正一次, 即按(2.8)式迭代一次得 y_{n+1} , 这个结果称校正值, 而这样建立的预测-校正系统通常称为改进的欧拉公式:

预测 $\tilde{y}_{n+1} = y_n + hf(x_n, y_n),$

校正 $y_{n+1} = y_n + \frac{h}{2}[f(x_n, y_n) + f(x_{n+1}, \tilde{y}_{n+1})]. \quad (2.13)$

或表为下列平均化形式

$$y_p = y_n + hf(x_n, y_n),$$

$$y_c = y_n + hf(x_{n+1}, y_p),$$

$$y_{n+1} = \frac{1}{2}(y_p + y_c).$$

例 2 用改进的欧拉方法求解初值问题(2.2).

解 改进的欧拉公式为

$$\begin{aligned}y_p &= y_n + h \cdot y_n - \frac{2x_n}{y_n}, \\y_c &= y_n + h \cdot y_p - \frac{2x_{n+1}}{y_p}, \\y_{n+1} &= \frac{1}{2}(y_p + y_c).\end{aligned}$$

仍取 $h=0.1$, 计算结果见表 9-2. 同例 1 中欧拉法的计算结果比较, 改进欧拉法明显改善了精度.

表 9-2 计算结果对比

x_n	y_n	$y(x_n)$	x_n	y_n	$y(x_n)$
0.1	1.0959	1.0954	0.6	1.4860	1.4832
0.2	1.1841	1.1832	0.7	1.5525	1.5492
0.3	1.2662	1.2649	0.8	1.6153	1.6165
0.4	1.3434	1.3416	0.9	1.6782	1.6733
0.5	1.4164	1.4142	1.0	1.7379	1.7321

9.3 龙格-库塔方法

9.3.1 显式龙格-库塔法的一般形式

上节给出了显式单步法的表达式(2.10), 其局部截断误差为(2.12), 对欧拉法 $T_{n+1} = O(h^2)$, 即方法为 $p=1$ 阶, 若用改进欧拉法(2.13), 它可表为

$$y_{n+1} = y_n + \frac{h}{2} [f(x_n, y_n) + f(x_n + h, y_n + hf(x_n, y_n))]. \quad (3.1)$$

此时增量函数

$$(x_n, y_n, h) = \frac{1}{2} [f(x_n, y_n) + f(x_n + h, y_n + hf(x_n, y_n))]. \quad (3.2)$$

它比欧拉法的 $(x_n, y_n, h) = f(x_n, y_n)$, 增加了计算一个右函数 f 的值, 可望 $p=2$. 若要使得到的公式阶数 p 更大, 就必须包含更多的 f 值. 实际上从方程(1.1)等价的积分形式(2.4), 即

$$y(x_{n+1}) - y(x_n) = \int_{x_n}^{x_{n+1}} f(x, y(x)) dx, \quad (3.3)$$

若要使公式阶数提高, 就必须使右端积分的数值求积公式精度提高, 它必然要增加求积节点, 为此可将(3.3)的右端用求积公式表示为

$$\int_{x_n}^{x_{n+1}} f(x, y(x)) dx = h \sum_{i=1}^r c_i f(x_n + i h, y(x_n + i h)).$$

一般来说, 点数 r 越多, 精度越高, 上式右端相当于增量函数 (x, y, h) , 为得到便于计算的显式方法, 可类似于改进欧拉法(3.1), (3.2), 将公式表示为

$$y_{n+1} = y_n + h (x_n, y_n, h), \quad (3.4)$$

其中

$$(x_n, y_n, h) = \sum_{i=1}^r c_i K_i, \quad (3.5)$$

$$K_1 = f(x_n, y_n),$$

$$K_i = f(x_n + i h, y_n + h \sum_{j=1}^{i-1} \mu_{ij} K_j) \quad i = 2, \dots, r,$$

这里 c_i, μ_{ij} 均为常数.(3.4)和(3.5)称为 r 级显式龙格-库塔(Runge-Kutta)法, 简称 R-K 方法.

当 $r=1$, $(x_n, y_n, h) = f(x_n, y_n)$ 时, 就是欧拉法, 此时方法的阶为 $p=1$. 当 $r=2$ 时, 改进欧拉法(3.1), (3.2)就是其中的一种, 下面将证明阶 $p=2$. 要使公式(3.4), (3.5)具有更高的阶 p , 就要增加点数 r . 下面我们只就 $r=2$ 推导 R-K 方法. 并给出 $r=3, 4$ 时

的常用公式, 其推导方法与 $r=2$ 时类似, 只是计算较复杂.

9.3.2 二阶显式 R-K 方法

对 $r=2$ 的 R-K 方法, 由(3.4), (3.5)可得到如下的计算公式

$$\begin{aligned} y_{n+1} &= y_n + h(a K_1 + c K_2), \\ K_1 &= f(x_n, y_n), \\ K_2 &= f(x_n + \frac{1}{2}h, y_n + \mu_1 h K_1). \end{aligned} \quad (3.6)$$

这里 a, c, μ_1 均为待定常数, 我们希望适当选取这些系数, 使公式阶数 p 尽量高. 根据局部截断误差定义, (3.6)的局部截断误差为

$$T_{n+1} = y(x_{n+1}) - y(x_n) - h[a f(x_n, y_n) + c f(x_n + \frac{1}{2}h, y_n + \mu_1 h f_n)], \quad (3.7)$$

这里 $y_n = y(x_n)$, $f_n = f(x_n, y_n)$. 为得到 T_{n+1} 的阶 p , 要将上式各项在 (x_n, y_n) 处做泰勒展开, 由于 $f(x, y)$ 是二元函数, 故要用到二元泰勒展开, 各项展开式为

$$y(x_{n+1}) = y_n + hy_n + \frac{h^2}{2}y_n + \frac{h^3}{3!}y_n + O(h^4),$$

其中

$$\begin{aligned} y_n &= f(x_n, y_n) = f_n, \\ y_n &= \frac{d}{dx}f(x_n, y(x_n)) = f_x(x_n, y_n) + f_y(x_n, y_n) \cdot f_n, \\ y_n &= f_{xx}(x_n, y_n) + 2f_n f_{xy}(x_n, y_n) + f_n^2 f_{yy}(x_n, y_n) \\ &\quad + f_y(x_n, y_n)[f_x(x_n, y_n) + f_n f_y(x_n, y_n)]; \end{aligned} \quad (3.8)$$

$f(x_n + \frac{1}{2}h, y_n + \mu_1 h f_n)$
 $= f_n + f_x(x_n, y_n) \frac{1}{2}h + f_y(x_n, y_n) \mu_1 h f_n + O(h^2)$. 将以上结果代入(3.7)则有

$$T_{n+1} = hf_n + \frac{h^2}{2}[f_x(x_n, y_n) + f_y(x_n, y_n) f_n]$$

$$\begin{aligned}
& - h[\alpha f_n + \alpha_1 (f_n + \frac{1}{2} f_x(x_n, y_n) h \\
& + \mu_1 f_y(x_n, y_n) f_n h) + O(h^3)] \\
& = (1 - \alpha - \alpha_1) f_n h + \frac{1}{2} - \alpha_1 \frac{1}{2} f_x(x_n, y_n) h^2 \\
& + \frac{1}{2} - \alpha \mu_1 f_y(x_n, y_n) f_n h^2 + O(h^3).
\end{aligned}$$

要使公式(3.6)具有 $p=2$ 阶, 必须使

$$1 - \alpha - \alpha_1 = 0, \quad \frac{1}{2} - \alpha_1 \frac{1}{2} = 0, \quad \frac{1}{2} - \alpha \mu_1 = 0. \quad (3.9)$$

即 $\alpha_1 = \frac{1}{2}, \quad \alpha \mu_1 = \frac{1}{2}, \quad \alpha + \alpha_1 = 1.$

(3.9)的解是不唯一的. 可令 $\alpha = a \neq 0$, 则得

$$\alpha = 1 - a, \quad \alpha_1 = \mu_1 = \frac{1}{2a}.$$

这样得到的公式称为二阶 R-K 方法, 如取 $a = 1/2$, 则 $\alpha = \alpha_1 = 1/2, \quad \mu_1 = 1$. 这就是改进欧拉法(3.1).

若取 $a = 1$, 则 $\alpha = 1, \alpha_1 = 0, \mu_1 = 1/2$. 得计算公式

$$\begin{aligned}
y_{n+1} &= y_n + hK_2, \\
K_1 &= f(x_n, y_n), \\
K_2 &= f(x_n + \frac{h}{2}, y_n + \frac{h}{2} K_1).
\end{aligned} \quad (3.10)$$

称为中点公式, 相当于数值积分的中矩形公式.(3.10)也可表示为

$$y_{n+1} = y_n + hf(x_n + \frac{h}{2}, y_n + \frac{h}{2} f(x_n, y_n)).$$

对 $r=2$ 的 R-K 公式(3.6)能否使局部误差提高到 $O(h^4)$? 为此需把 K_2 多展开一项, 从(3.8)的 y_n 看到展开式中 $f_y f_x + f_y f$ 的项是不能通过选择参数消掉的, 实际上要使 h^3 的项为零, 需增加 3 个方程, 要确定 4 个参数 α, α_1, μ_1 及 μ_2 , 这是不可能的. 故 $r=2$ 的显式 R-K 方法的阶只能是 $p=2$, 而不能得到三阶公式.

9.3.3 三阶与四阶显式 R-K 方法

要得到三阶显式 R-K 方法, 必须 $r=3$. 此时(3.4), (3.5)的公式表示为

$$\begin{aligned} y_{n+1} &= y_n + h(a K_1 + c_2 K_2 + c_3 K_3), \\ K_1 &= f(x_n, y_n), \\ K_2 &= f(x_n + \frac{1}{2}h, y_n + \mu_1 h K_1), \\ K_3 &= f(x_n + \frac{2}{3}h, y_n + \mu_1 h K_1 + \mu_2 h K_2). \end{aligned} \quad (3.11)$$

其中 a, c_2, c_3 及 μ_1, μ_2 均为待定参数, 公式(3.11)的局部截断误差为

$$T_{n+1} = y(x_{n+1}) - y(x_n) - h[a K_1 + c_2 K_2 + c_3 K_3].$$

只要将 K_2, K_3 按二元函数泰勒展开, 使 $T_{n+1} = O(h^4)$, 可得待定参数满足方程

$$\begin{aligned} a + c_2 + c_3 &= 1, \\ c_2 &= \mu_1, \\ c_3 &= \mu_1 + \mu_2, \\ c_2 + c_3 &= \frac{1}{2}, \\ c_2^2 + c_3^2 &= \frac{1}{3}, \\ c_2 \mu_2 &= \frac{1}{6}. \end{aligned} \quad (3.12)$$

这是 8 个未知数 6 个方程的方程组, 解也不是唯一的. 可以得到很多公式. 满足条件(3.12)的公式(3.11)统称为三阶 R-K 公式. 下面只给出其中一个常见的公式.

$$y_{n+1} = y_n + \frac{h}{6} (K_1 + 4K_2 + K_3),$$

$$K_1 = f(x_n, y_n),$$

$$K_2 = f(x_n + \frac{h}{2}, y_n + \frac{h}{2} K_1),$$

$$K_3 = f(x_n + h, y_n - hK_1 + 2hK_2).$$

此公式称为库塔三阶方法.

继续上述过程, 经过较复杂的数学演算, 可以导出各种四阶龙格-库塔公式, 下列经典公式是其中常用的一个:

$$y_{n+1} = y_n + \frac{h}{6} (K_1 + 2K_2 + 2K_3 + K_4),$$

$$K_1 = f(x_n, y_n),$$

$$K_2 = f(x_n + \frac{h}{2}, y_n + \frac{h}{2} K_1), \quad (3.13)$$

$$K_3 = f(x_n + \frac{h}{2}, y_n + \frac{h}{2} K_2),$$

$$K_4 = f(x_n + h, y_n + hK_3).$$

四阶龙格-库塔方法的每一步需要计算四次函数值 f , 可以证明其截断误差为 $O(h^5)$. 不过证明极其繁琐, 这里从略.

例 3 设取步长 $h=0.2$, 从 $x=0$ 直到 $x=1$ 用四阶龙格-库塔方法求解初值问题(2.2).

解 这里, 经典的四阶龙格-库塔公式(3.13)具有形式

$$y_{n+1} = y_n + \frac{h}{6} (K_1 + 2K_2 + 2K_3 + K_4),$$

$$K_1 = y_n - \frac{2x_n}{y_n},$$

$$K_2 = y_n + \frac{h}{2} K_1 - \frac{2x_n + h}{y_n + \frac{h}{2} K_1},$$

$$K_3 = y_n + \frac{h}{2} K_2 - \frac{2x_n + h}{y_n + \frac{h}{2} K_2},$$

$$K_4 = y_n + hK_3 - \frac{2(x_n + h)}{y_n + hK_3}.$$

右面列出计算结果 y_n , 表中 $y(x_n)$ 仍表示准确解 .

比较例 3 和例 2 的计算结果, 显然

以龙格-库塔方法的精度为高 .要注意, 虽然四阶龙格-库塔方法的计算量(每一步要 4 次计算函数 f)比改进的欧拉方法(它是一种二阶龙格-库塔方法, 每一步只要 2 次计算函数 f)大一倍, 但由于这里放大了步长($h = 0.2$), 表 9-3 和表 9-2 所耗费的计算量几乎相同 .这个例子又一次显示了选择算法的重要意义 .

然而值得指出的是, 龙格-库塔方法的推导基于泰勒展开方法, 因而它要求所求的解具有较好的光滑性质 .反之, 如果解的光滑性差, 那么, 使用四阶龙格-库塔方法求得的数值解, 其精度可能反而不如改进的欧拉方法 .实际计算时, 我们应当针对问题的具体特点选择合适的算法 .

表 9-3 计算结果

x_n	y_n	$y(x_n)$
0.2	1.1832	1.1832
0.4	1.3417	1.3416
0.6	1.4833	1.4832
0.8	1.6125	1.6125
1.0	1.7321	1.7321

9.3.4 变步长的龙格-库塔方法

单从每一步看, 步长越小, 截断误差就越小, 但随着步长的缩小, 在一定求解范围内所要完成的步数就增加了. 步数的增加不但引起计算量的增大, 而且可能导致舍入误差的严重积累. 因此同积分的数值计算一样, 微分方程的数值解法也有个选择步长的问题.

在选择步长时, 需要考虑两个问题:

1° 怎样衡量和检验计算结果的精度?

2° 如何依据所获得的精度处理步长?

我们考察经典的四阶龙格-库塔公式(3.13). 从节点 x_n 出发, 先以 h 为步长求出一个近似值, 记为 $y_{n+1}^{(h)}$, 由于公式的局部截断误差为 $O(h^5)$, 故有

$$y(x_{n+1}) - y_{n+1}^{(h)} = ch^5, \quad (3.14)$$

然后将步长折半, 即取 $\frac{h}{2}$ 为步长从 x_n 跨两步到 x_{n+1} , 再求得一个近似值 $y_{n+1}^{\frac{h}{2}}$, 每跨一步的截断误差是 $c \frac{h}{2}^5$, 因此有

$$y(x_{n+1}) - y_{n+1}^{\frac{h}{2}} = 2c \frac{h}{2}^5, \quad (3.15)$$

比较(3.14)式和(3.15)式我们看到, 步长折半后, 误差大约减少到 $\frac{1}{16}$, 即有

$$\frac{y(x_{n+1}) - y_{n+1}^{\frac{h}{2}}}{y(x_{n+1}) - y_{n+1}^{(h)}} = \frac{1}{16}.$$

由此易得下列事后估计式

$$y(x_{n+1}) - y_{n+1}^{\frac{h}{2}} = \frac{1}{15} y_{n+1}^{\frac{h}{2}} - y_{n+1}^{(h)}.$$

这样, 我们可以通过检查步长, 折半前后两次计算结果的偏差

$$= |y_{n+1}^{\frac{h}{2}} - y_{n+1}^{(h)}|$$

来判定所选的步长是否合适, 具体地说, 将区分以下两种情况处理:

1. 对于给定的精度, 如果 $\epsilon > \delta$, 我们反复将步长折半进行计算, 直至 $\epsilon < \delta$ 为止, 这时取最终得到的 $y_{n+1}^{\frac{h}{2}}$ 作为结果;

2. 如果 $\epsilon < \delta$, 我们将反复将步长加倍, 直到 $\epsilon > \delta$ 为止, 这时再将步长折半一次, 就得到所要的结果.

这种通过加倍或折半处理步长的方法称为变步长方法. 表面上看, 为了选择步长, 每一步的计算量增加了, 但总体考虑往往是合算的.

9.4 单步法的收敛性与稳定性

9.4.1 收敛性与相容性

数值解法的基本思想是, 通过某种离散化手段将微分方程(1.1)转化为差分方程, 如单步法(2.10), 即

$$y_{n+1} = y_n + h(x_n, y_n, h). \quad (4.1)$$

它在 x_n 处的解为 y_n , 而初值问题(1.1), (1.2)在 x_n 处的精确解为 $y(x_n)$, 记 $e_n = y(x_n) - y_n$ 称为整体截断误差. 收敛性就是讨论当

$x = x_n$ 固定且 $h = \frac{x_n - x_0}{n} \rightarrow 0$ 时 $e_n \rightarrow 0$ 的问题.

定义3 若一种数值方法(如单步法(4.1))对于固定的 $x_n = x_0 + nh$, 当 $h \rightarrow 0$ 时有 $y_n \rightarrow y(x_n)$, 其中 $y(x)$ 是(1.1), (1.2)的准确解, 则称该方法是收敛的.

显然数值方法收敛是指 $e_n = y(x_n) - y_n \rightarrow 0$, 对单步法(4.1)有下述收敛性定理:

定理1 假设单步法(4.1)具有 p 阶精度, 且增量函数 $f(x, y, h)$ 关于 y 满足利普希茨条件

$$| (x, y, h) - (\bar{x}, \bar{y}, h) | \leq L |y - \bar{y}|, \quad (4.2)$$

又设初值 y_0 是准确的, 即 $y_0 = y(x_0)$, 则其整体截断误差

$$|y(x_n) - y_n| = O(h^p). \quad (4.3)$$

证明 设以 \bar{y}_{n+1} 表示取 $y_n = y(x_n)$ 用公式(4.1)求得的结果, 即

$$\bar{y}_{n+1} = y(x_n) + h (x_n, y(x_n), h), \quad (4.4)$$

则 $|y(x_{n+1}) - \bar{y}_{n+1}|$ 为局部截断误差, 由于所给方法具有 p 阶精度, 按定义 2, 存在定数 C , 使

$$|y(x_{n+1}) - \bar{y}_{n+1}| \leq Ch^{p+1}.$$

又由式(4.4)与(4.1), 得

$$\begin{aligned} |\bar{y}_{n+1} - y_{n+1}| &= |y(x_n) - y_n| \\ &\quad + h |(x_n, y(x_n), h) - (x_n, y_n, h)|. \end{aligned}$$

利用假设条件(4.2), 有

$$|\bar{y}_{n+1} - y_{n+1}| \leq (1 + hL) |y(x_n) - y_n|,$$

从而有

$$\begin{aligned} |y(x_{n+1}) - y_{n+1}| &\leq |\bar{y}_{n+1} - y_{n+1}| + |y(x_{n+1}) - \bar{y}_{n+1}| \\ &\leq (1 + hL) |y(x_n) - y_n| + Ch^{p+1}. \end{aligned}$$

即对整体截断误差 $e_n = y(x_n) - y_n$ 成立下列递推关系式

$$|e_{n+1}| \leq (1 + hL) |e_n| + Ch^{p+1}, \quad (4.5)$$

据此不等式反复递推, 可得

$$|e_n| \leq (1 + hL)^n |e_0| + \frac{Ch^p}{L} [(1 + hL)^n - 1]. \quad (4.6)$$

再注意到当 $x_n - x_0 = nh = T$ 时

$$(1 + hL)^n = (e^{hL})^n = e^{TL},$$

最终得下列估计式

[注] 对于任意实数 x , 有 $1 + x \leq e^x$, 而当 $x \geq 1$ 时, 成立 $0 < (1 + x)^n \leq e^{nx}$.

$$|e_n| / |e_0| e^{TL} + \frac{Ch^p}{L} (e^{TL} - 1). \quad (4.7)$$

由此可以断定, 如果初值是准确的, 即 $e_0 = 0$, 则(4.3)式成立. 定理证毕.

依据这一定理, 判断单步法(4.1)的收敛性, 归结为验证增量函数能否满足利普希茨条件(4.2).

对于欧拉方法, 由于其增量函数就是 $f(x, y)$, 故当 $f(x, y)$ 关于 y 满足利普希茨条件时它是收敛的.

再考察改进的欧拉方法, 其增量函数已由(3.2)式给出, 这时有

$$\begin{aligned} & |(x, y, h) - (x, \tilde{y}, h)| = \frac{1}{2} [|f(x, y) - f(x, \tilde{y})| \\ & + |f(x + h, y + hf(x, y)) - f(x + h, \tilde{y} + hf(x, \tilde{y}))|]. \end{aligned}$$

假设 $f(x, y)$ 关于 y 满足利普希茨条件, 记利普希茨常数为 L , 则由上式推得

$$|(x, y, h) - (x, \tilde{y}, h)| \leq L \left(1 + \frac{h}{2} L\right) |y - \tilde{y}|.$$

设限定 $h = h_0$ (h_0 为定数), 上式表明 关于 y 的利普希茨常数

$$L' = L \left(1 + \frac{h_0}{2} L\right),$$

因此改进的欧拉方法也是收敛的.

类似地, 不难验证其他龙格-库塔方法的收敛性.

定理 1 表明 $p=1$ 时单步法收敛, 并且当 $y(x)$ 是初值问题(1.1), (1.2)的解, (4.1)具有 p 阶精度时, 则有展开式

$$\begin{aligned} T_{n+1} &= y(x + h) - y(x) - h (x, y(x), h) \\ &= y(x)h + \frac{y'(x)}{2} h^2 + \dots \\ &\quad - h[(x, y(x), 0) + {}_x(x, y(x), 0)h + \dots] \\ &= h[y(x) - (x, y(x), 0)] + O(h^2). \end{aligned}$$

所以 $p=1$ 的充要条件是 $y(x) - (x, y(x), 0) = 0$, 而 $y(x) = f(x, y(x))$, 于是可给出如下定义:

定义 4 若单步法(4.1)的增量函数 满足

$$(x, y, 0) = f(x, y),$$

则称单步法(4.1)与初值问题(1.1),(1.2)相容.

以上讨论表明 p 阶方法(4.1)当 $p=1$ 时与(1.1),(1.2)相容, 反之相容方法至少是 1 阶的.

于是由定理 1 可知方法(4.1)收敛的充分必要条件是此方法是相容的.

9.4.2 绝对稳定性与绝对稳定域

前面关于收敛性的讨论有个前提, 必须假定数值方法本身的计算是准确的. 实际情形并不是这样, 差分方程的求解还会有计算误差, 譬如由于数字舍入而引起的小扰动. 这类小扰动在传播过程中会不会恶性增长, 以至于“淹没”了差分方程的“真解”呢? 这就是差分方法的稳定性问题. 在实际计算时, 我们希望某一步产生的扰动值, 在后面的计算中能够被控制, 甚至是逐步衰减的.

定义 5 若一种数值方法在节点值 y_n 上大小为 的扰动, 于以后各节点值 y_m ($m > n$) 上产生的偏差均不超过 , 则称该方法是稳定的.

下面先以欧拉法为例考察计算稳定性.

例 4 考察初值问题

$$y' = -100y,$$

$$y(0) = 1.$$

其准确解 $y(x) = e^{-100x}$ 是一个按指数曲线衰减得很快的函数, 如图 9-3 所示.

用欧拉法解方程 $y' = -100y$ 得

$$y_{n+1} = (1 - 100h)y_n.$$

若取 $h = 0.025$, 则欧拉公式的具体形式为

$$y_{n+1} = -1.5 y_n,$$

计算结果列于表 9-4 的第 2 列. 我们看到, 欧拉方法的解 y_n (图 9-3 中用 \times 号标出) 在准确值 $y(x_n)$ 的上下波动, 计算过程明显地不稳定. 但若取 $h = 0.005$, $y_{n+1} = 0.5 y_n$ 则计算过程稳定.

再考察后退的欧拉方法, 取 $h = 0.025$ 时计算公式为

$$y_{n+1} = \frac{1}{3.5} y_n.$$

计算结果列于表 9-4 的第 3 列(图 9-3 中标以 · 号), 这时计算过程是稳定的.

表 9-4 计算结果对比

节 点	欧拉方法	后退欧拉方法
0.025	-1.5	0.2857
0.050	2.25	0.0816
0.075	-3.375	0.0233
0.100	5.0625	0.0067

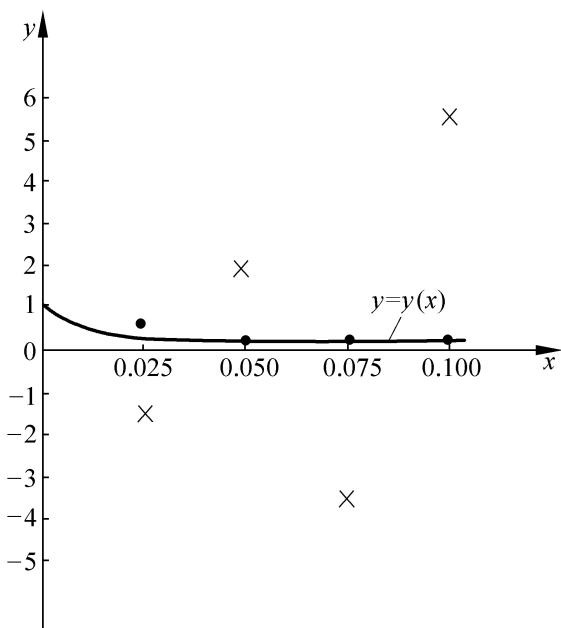


图 9-3

例题表明稳定性不但与方法有关, 也与步长 h 的大小有关, 当然也与方程中的 $f(x, y)$ 有关. 为了只考察数值方法本身. 通常只检验将数值方法用于解模型方程的稳定性, 模型方程为

$$y' = -y, \quad (4.8)$$

其中 y 为复数, 这个方程分析较简单. 对一般方程可以通过局部线性化化为这种形式, 例如在 (\bar{x}, \bar{y}) 的邻域, 可展开为

$$\begin{aligned} y &= f(x, y) \\ &= f(\bar{x}, \bar{y}) + f_x(\bar{x}, \bar{y})(x - \bar{x}) + f_y(\bar{x}, \bar{y})(y - \bar{y}) + \dots, \end{aligned}$$

略去高阶项, 再做变换即可得到 $u = u$ 的形式. 对于 m 个方程的方程组, 可线性化为 $\mathbf{y} = \mathbf{Ay}$, 这里 \mathbf{A} 为 $m \times m$ 的雅可比矩阵

$\frac{f_i}{y_j}$. 若 \mathbf{A} 有 m 个特征值 $\lambda_1, \dots, \lambda_m$, 其中 λ_i 可能是复数, 所以,

为了使模型方程结果能推广到方程组, 方程(4.8)中 y 为复数. 为保证微分方程本身的稳定性, 还应假定 $\operatorname{Re}(\lambda_i) < 0$.

下面先研究欧拉方法的稳定性. 模型方程 $y' = y$ 的欧拉公式为

$$y_{n+1} = (1 + h)y_n. \quad (4.9)$$

设在节点值 y_n 上有一扰动值 ϵ_n , 它的传播使节点值 y_{n+1} 产生大小为 ϵ_{n+1} 的扰动值, 假设用 $y_n^* = y_n + \epsilon_n$ 按欧拉公式得出 $y_{n+1}^* = y_{n+1} + \epsilon_{n+1}$ 的计算过程不再有新的误差, 则扰动值满足

$$\epsilon_{n+1} = (1 + h)\epsilon_n.$$

可见扰动值满足原来的差分方程(4.9). 这样, 如果差分方程的解是不增长的, 即有

$$|y_{n+1}| \leq |y_n|,$$

则它就是稳定的. 这一论断对于下面将要研究的其他方法同样适用.

显然, 为要保证差分方程(4.9)的解是不增长的, 只要选取 h 充分小, 使

$$|1 + h| \leq 1. \quad (4.10)$$

在 $\mu = h$ 的复平面上, 这是以 $(-1, 0)$ 为圆心, 1 为半径的单位圆域. 称为欧拉法的绝对稳定域, 一般情形可由下面定义.

定义 6 单步法(4.1)用于解模型方程(4.8), 若得到的解

$y_{n+1} = E(h)y_n$, 满足 $|E(h)| < 1$, 则称方法(4.1)是绝对稳定的. 在 $\mu = h$ 的平面上, 使 $|E(h)| < 1$ 的变量围成的区域, 称为绝对稳定域, 它与实轴的交称为绝对稳定区间.

对欧拉法 $E(h) = 1 + h$, 其绝对稳定域已由(4.10)给出, 绝对稳定区间为 $-2 < h < 0$, 在例 5 中 $= -100$, $-2 < -100h < 0$, 即 $0 < h < 2/100 = 0.02$ 为绝对稳定区间, 例 4 中取 $h = 0.025$ 故它是不稳定的, 当取 $h = 0.005$ 时它是稳定的.

对二阶 R-K 方法, 解模型方程(4.1)可得到

$$y_{n+1} = 1 + h + \frac{(h)^2}{2} y_n,$$

故 $E(h) = 1 + h + \frac{(h)^2}{2}$.

绝对稳定域由 $\left|1 + h + \frac{(h)^2}{2}\right| < 1$ 得到, 于是可得绝对稳定区间为 $-2 < h < 0$, 即 $0 < h < -2$. 类似可得三阶及四阶的 R-K 方法的 $E(h)$ 分别为

$$E(h) = 1 + h + \frac{(h)^2}{2!} + \frac{(h)^3}{3!},$$

$$E(h) = 1 + h + \frac{(h)^2}{2!} + \frac{(h)^3}{3!} + \frac{(h)^4}{4!}.$$

由 $|E(h)| < 1$ 可得到相应的绝对稳定域. 当 h 为实数时则得绝对稳定区间. 它们分别为

三阶显式 R-K 方法: $-2.51 < h < 0$, 即 $0 < h < -2.51$.

四阶显式 R-K 方法: $-2.78 < h < 0$, 即 $0 < h < -2.78$.

从以上讨论可知显式的 R-K 方法的绝对稳定域均为有限域, 都对步长 h 有限制. 如果 h 不在所给的绝对稳定区间内, 方法就不稳定.

例 5 $y' = -20y$ ($0 < x < 1$), $y(0) = 1$, 分别取 $h = 0.1$ 及 $h = 0.2$ 用经典的四阶 R-K 方法(3.13)计算.

解 本例 $\mu = -20$, h 分别为 -2 及 -4, 前者在绝对稳定区间内, 后者则不在, 用四阶 R-K 方法计算其误差见下表:

x_n	0.2	0.4	0.6	0.8	1.0
$h=0.1$	0.93×10^{-1}	0.12×10^{-1}	0.14×10^{-2}	0.15×10^{-3}	0.17×10^{-4}
$h=0.2$	4.98	25.0	125.0	625.0	3125.0

从以上结果看到, 如果步长 h 不满足绝对稳定条件, 误差增长很快.

对隐式单步法, 可以同样讨论方法的绝对稳定性, 例如对后退欧拉法, 用它解模型方程可得

$$y_{n+1} = \frac{1}{1 - h} y_n,$$

故 $E(h) = \frac{1}{1 - h}$.

由 $|E(h)| = \left| \frac{1}{1 - h} \right| < 1$ 可得绝对稳定域为 $|1 - h| > 1$, 它是以 $(1, 0)$ 为圆心, 1 为半径的单位圆外部, 故绝对稳定区间为 $-1 < h < 0$. 当 $h < 0$ 时, 则 $0 < |h| < 1$, 即对任何步长均为稳定的.

对隐式梯形法, 它用于解模型方程(4.8)得

$$y_{n+1} = \frac{1 + \frac{h}{2}}{1 - \frac{h}{2}} y_n,$$

故 $E(h) = \frac{1 + h/2}{1 - h/2}$.

对 $\operatorname{Re}(h) < 0$ 有 $|E(h)| = \left| \frac{1 + \frac{h}{2}}{1 - \frac{h}{2}} \right| < 1$, 故绝对稳定域为 $\mu = h$

的左半平面, 绝对稳定区间为 $- \frac{1}{2} < h < 0$, 即 $0 < h < \frac{1}{2}$ 时梯形法均是稳定的 .

9.5 线性多步法

在逐步推进的求解过程中, 计算 y_{n+1} 之前事实上已经求出了一系列的近似值 y_0, y_1, \dots, y_n , 如果充分利用前面多步的信息来预测 y_{n+1} , 则可以期望会获得较高的精度 . 这就是构造所谓线性多步法的基本思想 .

构造多步法的主要途径是基于数值积分方法和基于泰勒展开方法, 前者可直接由方程(1.1)两端积分后利用插值求积公式得到 . 本节主要介绍基于泰勒展开的构造方法 .

9.5.1 线性多步法的一般公式

如果计算 y_{n+k} 时, 除用 y_{n+k-1} 的值, 还用到 y_{n+i} ($i = 0, 1, \dots, k-2$) 的值, 则称此方法为线性多步法 . 一般的线性多步法公式可表示为

$$y_{n+k} = \sum_{i=0}^{k-1} c_i y_{n+i} + h \sum_{i=0}^k f_{n+i}, \quad (5.1)$$

其中 y_{n+i} 为 $y(x_{n+i})$ 的近似, $f_{n+i} = f(x_{n+i}, y_{n+i})$, $x_{n+i} = x_0 + ih$, c_i 为常数, c_0 及 c_k 不全为零, 则称(5.1)为线性 k 步法, 计算时需先给出前面 k 个近似值 y_0, y_1, \dots, y_{k-1} , 再由(5.1)逐次求出 y_k, y_{k+1}, \dots . 如果 $c_k = 0$, 称(5.1)为显式 k 步法, 这时 y_{n+k} 可直接由(5.1)算出; 如果 $c_k \neq 0$, 则(5.1)称为隐式 k 步法, 求解时与梯形法(2.7)相同, 要用迭代法方可算出 y_{n+k} . (5.1)中系数 c_i 及 f_i 可根据方法的局部截断误差及阶确定, 其定义为:

定义 7 设 $y(x)$ 是初值问题(1.1), (1.2)的准确解, 线性多步法(5.1)在 x_{n+k} 上的局部截断误差为

$$\begin{aligned}
 T_{n+k} &= L[y(x_n); h] \\
 &= y(x_{n+k}) - \sum_{i=0}^{k-1} y(x_{n+i}) - h \sum_{i=0}^k y(x_{n+i}). \quad (5.2)
 \end{aligned}$$

若 $T_{n+k} = O(h^{p+1})$, 则称方法(5.1)是 p 阶的, $p=1$ 则称方法(5.1)与方程(1.1)是相容的.

由定义7, 对 T_{n+k} 在 x_n 处做泰勒展开, 由于

$$\begin{aligned}
 y(x_n + ih) &= y(x_n) + ihy(x_n) + \frac{(ih)^2}{2!} y''(x_n) \\
 &\quad + \frac{(ih)^3}{3!} y'''(x_n) + \dots, \\
 y(x_n + ih) &= y(x_n) + ihy(x_n) + \frac{(ih)^2}{2!} y''(x_n) + \dots.
 \end{aligned}$$

代入(5.2)得

$$\begin{aligned}
 T_{n+k} &= a_0 y(x_n) + a_1 hy(x_n) + a_2 h^2 y''(x_n) \\
 &\quad + \dots + c_p h^p y^{(p)}(x_n) + \dots, \quad (5.3)
 \end{aligned}$$

其中

$$\begin{aligned}
 a_0 &= 1 - (0 + \dots + k-1), \\
 a_1 &= k - [1 + 2_2 + \dots + (k-1)_{k-1}] - (0 + 1 + \dots + k), \\
 c_q &= \frac{1}{q!} [k^q - (1 + 2^q_2 + \dots + (k-1)^q_{k-1}) \\
 &\quad - \frac{1}{(q-1)!} [1 + 2^{q-1}_2 + \dots + k^{q-1}_{k-1}]] \\
 q &= 2, 3, \dots. \quad (5.4)
 \end{aligned}$$

若在公式(5.1)中选择系数 a_i 及 c_i , 使它满足

$$a_0 = a_1 = \dots = c_p = 0, \quad c_{p+1} \neq 0.$$

由定义可知此时所构造的多步法是 p 阶的, 且

$$T_{n+k} = c_{p+1} h^{p+1} y^{(p+1)}(x_n) + O(h^{p+2}). \quad (5.5)$$

称右端第一项为局部截断误差主项, c_{p+1} 称为误差常数.

根据相容性定义, $p=1$, 即 $a_0 = a_1 = 0$, 由(5.4)得

$$\begin{aligned} 0 &+ \dots + k-1 = 1, \\ i &+ \dots + i = k. \end{aligned} \tag{5.6}$$

$i=1$ $i=0$

故方法(5.1)与微分方程(1.1)相容的充分必要条件是(5.6)成立.

显然,当 $k=1$ 时,若 $\beta_1=0$, 则由(5.6)可求得

$$_0 = 1, \quad _0 = 1.$$

此时公式(5.1)为

$$y_{n+1} = y_n + hf_n,$$

即为欧拉法. 从(5.4)可求得 $c_2 = 1/2 - 0$, 故方法为 1 阶精度, 且局部截断误差为

$$T_{n+1} = \frac{1}{2} h^2 y(x_n) + O(h^3),$$

这和第 2 节给出的定义及结果是一致的.

对 $k=1$, 若 $\alpha_1 = 0$, 此时方法为隐式公式, 为了确定系数 $\alpha_0, \alpha_1, \alpha_2$, 可由 $a_0 = a_1 = a_2 = 0$ 解得 $\alpha_0 = 1, \alpha_1 = \alpha_2 = 1/2$. 于是得到公式

$$y_{n+1} = y_n + \frac{h}{2} (f_n + f_{n+1}),$$

即为梯形法.由(5.4)可求得 $c_3 = -1/12$,故 $p=2$,所以梯形法是二阶方法,其局部截断误差主项是 $-h^3 y'(x_n)/12$,这与第2节中的讨论也是一致的.

对 $k=2$ 的多步法公式都可利用(5.4)确定系数 α_i, β_i , 并由(5.5)给出局部截断误差, 下面只就若干常用的多步法导出具体公式.

9.5.2 阿当姆斯显式与隐式公式

考慮形如

$$y_{n+k} = y_{n+k-1} + h \sum_{i=0}^k f_{n+i} \quad (5.7)$$

的 k 步法，称为阿当姆斯(Adams)方法。 $k = 0$ 为显式方法， $k > 0$ 为隐式方法。

为隐式方法,通常称为阿当姆斯显式与隐式公式,也称 Adams-Bashforth 公式与 Adams-Moulton 公式.这类公式可直接由方程(1.1)两端积分(从 x_{n+k-1} 到 x_{n+k} 积分)求得.下面可利用(5.4)由 $a_0 = \dots = a_p = 0$ 推出,对比(5.7)与(5.1)可知此时系数 $a_0 = a_1 = \dots = a_{k-2} = 0$, $a_{k-1} = 1$,显然 $a_k = 0$ 成立,下面只需确定系数 a_0, a_1, \dots, a_k ,故可令 $a_0 = \dots = a_{k+1} = 0$,则可求得 a_0, a_1, \dots, a_k (若 $a_k = 0$,则令 $a_0 = \dots = a_k = 0$ 来求得 a_0, a_1, \dots, a_{k-1}).下面以 $k=3$ 为例,由 $a_0 = a_1 = a_2 = a_3 = 0$,根据(5.4)得

$$\begin{aligned} a_0 + a_1 + a_2 + a_3 &= 1, \\ 2(a_1 + 2a_2 + 3a_3) &= 5, \\ 3(a_1 + 4a_2 + 9a_3) &= 19, \\ 4(a_1 + 8a_2 + 27a_3) &= 65. \end{aligned}$$

若 $a_3 = 0$,则由前三个方程解得

$$a_0 = \frac{5}{12}, \quad a_1 = -\frac{16}{12}, \quad a_2 = \frac{23}{12},$$

得到 $k=3$ 的阿当姆斯显式公式是

$$y_{n+3} = y_{n+2} + \frac{h}{12}(23f_{n+2} - 16f_{n+1} + 5f_n). \quad (5.8)$$

由(5.4)求得 $a_4 = 3/8$,所以(5.8)是三阶方法,局部截断误差是

$$T_{n+3} = \frac{3}{8}h^4 y^{(4)}(x_n) + O(h^5).$$

若 $a_3 \neq 0$,则可解得

$$a_0 = \frac{1}{24}, \quad a_1 = -\frac{5}{24}, \quad a_2 = \frac{19}{24}, \quad a_3 = \frac{3}{8}.$$

于是得 $k=3$ 的阿当姆斯隐式公式为

$$y_{n+3} = y_{n+2} + \frac{h}{24}(9f_{n+3} + 19f_{n+2} - 5f_{n+1} + f_n), \quad (5.9)$$

它是四阶方法,局部截断误差是

$$T_{n+3} = -\frac{19}{720}h^5 y^{(5)}(x_n) + O(h^6). \quad (5.10)$$

用类似的方法可求得阿当姆斯显式方法和隐式方法的公式, 表 9-5 及表 9-6 分别列出了 $k=1, 2, 3, 4$ 时的阿当姆斯显式公式与阿当姆斯隐式公式, 其中 k 为步数, p 为方法的阶, c_{p+1} 为误差常数.

表 9-5 阿当姆斯显式公式

k	p	公 式	c_{p+1}
1	1	$y_{n+1} = y_n + hf_n$	$\frac{1}{2}$
2	2	$y_{n+2} = y_{n+1} + \frac{h}{2}(3f_{n+1} - f_n)$	$\frac{5}{12}$
3	3	$y_{n+3} = y_{n+2} + \frac{h}{12}(23f_{n+2} - 16f_{n+1} + 5f_n)$	$\frac{3}{8}$
4	4	$y_{n+4} = y_{n+3} + \frac{h}{24}(55f_{n+3} - 59f_{n+2} + 37f_{n+1} - 9f_n)$	$\frac{251}{720}$

表 9-6 阿当姆斯隐式公式

k	p	公 式	c_{p+1}
1	2	$y_{n+1} = y_n + \frac{h}{2}(f_{n+1} + f_n)$	$-\frac{1}{12}$
2	3	$y_{n+2} = y_{n+1} + \frac{h}{12}(5f_{n+2} + 8f_{n+1} - f_n)$	$-\frac{1}{24}$
3	4	$y_{n+3} = y_{n+2} + \frac{h}{24}(9f_{n+3} + 19f_{n+2} - 5f_{n+1} + f_n)$	$-\frac{19}{720}$
4	5	$y_{n+4} = y_{n+3} + \frac{h}{720}(251f_{n+4} + 646f_{n+3} - 264f_{n+2} + 106f_{n+1} - 19f_n)$	$-\frac{3}{160}$

例 6 用四阶阿当姆斯显式和隐式方法解初值问题

$$y' = -y + x + 1, \quad y(0) = 1.$$

取步长 $h=0.1$.

解 本题 $f_n = -y_n + x_n + 1$, $x_n = nh = 0.1n$. 从四阶阿当姆斯

显式公式得到

$$\begin{aligned}
 y_{n+4} &= y_{n+3} + \frac{h}{24}(55 f_{n+3} - 59 f_{n+2} + 37 f_{n+1} - 9 f_n) \\
 &= \frac{1}{24}(18.5 y_{n+3} + 5.9 y_{n+2} - 3.7 y_{n+1} + 0.9 y_n \\
 &\quad + 0.24n + 3.24).
 \end{aligned}$$

对于四阶阿当姆斯隐式公式得到

$$\begin{aligned}
 y_{n+3} &= y_{n+2} + \frac{h}{24}(9 f_{n+3} + 19 f_{n+2} - 5 f_{n+1} + f_n) \\
 &= \frac{1}{24}(-0.9 y_{n+3} + 22.1 y_{n+2} + 0.5 y_{n+1} - 0.1 y_n + 0.24n + 3).
 \end{aligned}$$

由此可直接解出 y_{n+3} 而不用迭代, 得到

$$y_{n+3} = \frac{1}{24.9}(22.1 y_{n+2} + 0.5 y_{n+1} - 0.1 y_n + 0.24n + 3).$$

计算结果见表 9-7, 其中显式方法中的 y_0, y_1, y_2, y_3 及隐式方法中的 y_0, y_1, y_2 均用准确解 $y(x) = e^{-x} + x$ 计算得到, 对一般方程, 可用四阶 R-K 方法计算初始近似.

表 9-7 计算结果

x_n	精确解 $y(x_n)$ $= e^{-x_n} + x_n$	阿当姆斯显式方法		阿当姆斯隐式方法	
		y_n	$ y(x_n) - y_n $	y_n	$ y(x_n) - y_n $
0.3	1.040 818 22			1.040 818 01	2.1×10^{-7}
0.4	1.070 320 05	1.070 322 92	2.87×10^{-6}	1.070 319 66	3.9×10^{-7}
0.5	1.106 530 66	1.106 535 48	4.82×10^{-6}	1.106 530 14	5.2×10^{-7}
0.6	1.148 811 64	1.148 818 41	6.77×10^{-6}	1.148 811 01	6.3×10^{-7}
0.7	1.196 585 30	1.196 593 40	8.10×10^{-6}	1.196 584 59	7.1×10^{-7}
0.8	1.249 328 96	1.249 338 16	9.20×10^{-6}	1.249 328 19	7.7×10^{-7}
0.9	1.306 569 66	1.306 579 62	9.96×10^{-6}	1.306 568 84	8.2×10^{-7}
1.0	1.367 879 44	1.367 889 96	1.05×10^{-5}	1.367 878 59	8.5×10^{-7}

从以上例子看到同阶的阿当姆斯方法, 隐式方法要比显式方法误差小, 这可以从两种方法的局部截断误差主项 $c_{p+1} h^{p+1} y^{(p+1)}(x_n)$ 的系数大小得到解释, 这里 c_{p+1} 分别为 $25/720$ 及 $-19/720$.

9.5.3 米尔尼方法与辛普森方法

考虑与(5.7)不同的另一个 $k=4$ 的显式公式

$$y_{n+4} = y_n + h(-_3 f_{n+3} + _2 f_{n+2} + _1 f_{n+1} + _0 f_n),$$

其中 $_0, _1, _2, _3$ 为待定常数, 可根据使公式的阶尽可能高这一条件来确定其数值. 由(5.4)可知 $a_0 = 0$, 再令 $a_1 = a_2 = a_3 = a_4 = 0$ 得到

$$_0 + _1 + _2 + _3 = 4,$$

$$2(_1 + 2_2 + 3_3) = 16,$$

$$3(_1 + 4_2 + 9_3) = 64,$$

$$4(_1 + 8_2 + 27_3) = 256.$$

解此方程组得

$$_3 = \frac{8}{3}, \quad _2 = -\frac{4}{3}, \quad _1 = \frac{8}{3}, \quad _0 = 0.$$

于是得到四步显式公式

$$y_{n+4} = y_n + \frac{4h}{3}(2f_{n+3} - f_{n+2} + 2f_{n+1}), \quad (5.11)$$

称为米尔尼(Milne)方法. 由于 $c_5 = 14/45$, 故方法为 4 阶, 其局部截断误差为

$$T_{n+4} = \frac{14}{45}h^5 y^{(5)}(x_n) + O(h^6). \quad (5.12)$$

米尔尼方法也可以通过方程(1.1)两端积分

$$y(x_{n+4}) - y(x_n) = \int_{x_n}^{x_{n+4}} f(x, y(x)) dx$$

得到. 若将方程(1.1)从 x_n 到 x_{n+2} 积分, 可得

$$y(x_{n+2}) - y(x_n) = \int_{x_n}^{x_{n+2}} f(x, y(x)) dx.$$

右端积分通过辛普森求积公式就有

$$y_{n+2} = y_n + \frac{h}{3} (f_n + 4f_{n+1} + f_{n+2}). \quad (5.13)$$

称为辛普森方法. 它是隐式二步四阶方法, 其局部截断误差为

$$T_{n+2} = -\frac{h^5}{90} y^{(5)}(x_n) + O(h^6). \quad (5.14)$$

9.5.4 汉明方法

辛普森公式是二步方法中阶数最高的, 但它的稳定性较差, 为了改善稳定性, 我们考察另一类三步法公式

$$y_{n+3} = c_0 y_n + c_1 y_{n+1} + c_2 y_{n+2} + h(c_1 f_{n+1} + c_2 f_{n+2} + c_3 f_{n+3}),$$

其中系数 c_0, c_1, c_2 及 c_1, c_2, c_3 为常数, 如果希望导出的公式是四阶的, 则系数中至少有一个自由参数. 若取 $c_1 = 1$, 则可得到辛普森公式. 若取 $c_1 = 0$, 仍利用泰勒展开, 由(5.4), 令 $c_0 = c_1 = c_2 = c_3 = c_4 = 0$, 则可得到

$$c_0 + c_2 = 1,$$

$$2c_2 + c_1 + c_2 + c_3 = 3,$$

$$4c_2 + 2(c_1 + 2c_2 + 3c_3) = 9,$$

$$8c_2 + 3(c_1 + 4c_2 + 9c_3) = 27,$$

$$16c_2 + 4(c_1 + 8c_2 + 27c_3) = 81.$$

解此方程组得

$$c_0 = -\frac{1}{8}, \quad c_2 = \frac{9}{8}, \quad c_1 = -\frac{3}{8}, \quad c_2 = \frac{6}{8}, \quad c_3 = \frac{3}{8}.$$

于是有

$$y_{n+3} = \frac{1}{8}(9y_{n+2} - y_n) + \frac{3h}{8}(f_{n+3} + 2f_{n+2} - f_{n+1}), \quad (5.15)$$

称为汉明(Hamming)方法. 由于 $c = -1/40$, 故方法是四阶的, 且

局部截断误差为

$$T_{n+3} = -\frac{h^5}{40} y^{(5)}(x_n) + O(h^6). \quad (5.16)$$

9.5.5 预测-校正方法

对于隐式的线性多步法, 计算时要进行迭代, 计算量较大. 为了避免进行迭代, 通常采用显式公式给出 y_{n+k} 的一个初始近似, 记为 $y_{n+k}^{(0)}$, 称为预测(predictor), 接着计算 f_{n+k} 的值(evaluation), 再用隐式公式计算 y_{n+k} , 称为校正(corrector). 例如在(2.13)中用欧拉法做预测, 再用梯形法校正, 得到改进欧拉法, 它就是一个二阶预测-校正方法. 一般情况下, 预测公式与校正公式都取同阶的显式方法与隐式方法相匹配. 例如用四阶的阿当姆斯显式方法做预测, 再用四阶阿当姆斯隐式公式做校正, 得到以下格式:

$$\text{预测 P: } y_{n+4}^p = y_{n+3} + \frac{h}{24}(55f_{n+3} - 59f_{n+2} + 37f_{n+1} - 9f_n),$$

$$\text{求值 E: } f_{n+4}^p = f(x_{n+4}, y_{n+4}^p),$$

$$\text{校正 C: } y_{n+4} = y_{n+3} + \frac{h}{24}(9f_{n+4}^p + 19f_{n+3} - 5f_{n+2} + f_{n+1}),$$

$$\text{求值 E: } f_{n+4} = f(x_{n+4}, y_{n+4}).$$

此公式称为阿当姆斯四阶预测-校正格式(PECE).

依据四阶阿当姆斯公式的截断误差, 对于 PECE 的预测步 P 有

$$y(x_{n+4}) - y_{n+4}^p = \frac{251}{720} h^5 y^{(5)}(x_n),$$

对校正步 C 有

$$y(x_{n+4}) - y_{n+4} = -\frac{19}{720} h^5 y^{(5)}(x_n).$$

两式相减得

$$h^5 y^{(5)}(x_n) = -\frac{720}{270} (y_{n+4}^p - y_{n+4}),$$

于是有下列事后误差估计

$$y(x_{n+4}) - y_{n+4}^p - \frac{251}{270}(y_{n+4}^p - y_{n+4}),$$

$$y(x_{n+4}) - y_{n+4} - \frac{19}{270}(y_{n+4}^p - y_{n+4}).$$

容易看出

$$\begin{aligned} y_{n+4}^{pm} &= y_{n+4}^p + \frac{251}{270}(y_{n+4} - y_{n+4}^p), \\ \text{约} &= y_{n+4} - \frac{19}{270}(y_{n+4}^p - y_{n+4}). \end{aligned} \quad (5.17)$$

比 y_{n+4}^p , y_{n+4} 更好. 但在 y_{n+4}^{pm} 的表达式中 y_{n+4} 是未知的, 因此计算时用上一步代替, 从而构造一种修正预测-校正格式(PMECME):

$$P: y_{n+4}^p = y_{n+3} + \frac{h}{24}(55f_{n+3} - 59f_{n+2} + 37f_{n+1} - 9f_n),$$

$$M: y_{n+4}^{pm} = y_{n+4}^p + \frac{251}{270}(y_{n+3}^c - y_{n+3}^p),$$

$$E: f_{n+4}^{pm} = f(x_{n+4}, y_{n+4}^{pm}),$$

$$C: y_{n+4}^c = y_{n+3} + \frac{h}{24}(9f_{n+4}^{pm} + 19f_{n+3} - 5f_{n+2} + f_{n+1}),$$

$$M: y_{n+4} = y_{n+4}^c - \frac{19}{270}(y_{n+4}^c - y_{n+4}^p),$$

$$E: f_{n+4} = f(x_{n+4}, y_{n+4}).$$

注意: 在 PMECME 格式中已将(5.17)的 y_{n+4} 及 y_{n+4} 分别改为 y_{n+4}^c 及 y_{n+4} .

利用米尔尼公式(5.11)和汉明公式(5.15)相匹配, 并利用截断误差(5.12), (5.16)改进计算结果, 可类似地建立四阶修正米尔尼-汉明预测-校正格式(PMECME):

$$P: y_{n+4}^p = y_n + \frac{4}{3}h(2f_{n+3} - f_{n+2} + 2f_{n+1}),$$

$$M: y_{n+4}^{pm} = y_{n+4}^p + \frac{112}{121}(y_{n+3}^c - y_{n+3}^p),$$

$$E: f_{n+4}^{p,m} = f(x_{n+4}, y_{n+4}^{p,m}),$$

$$C: y_{n+4}^c = \frac{1}{8}(9y_{n+3} - y_{n+1}) + \frac{3}{8}h(f_{n+4}^{p,m} + 2f_{n+3} - f_{n+2}),$$

$$M: y_{n+4} = y_{n+4}^c - \frac{9}{121}(y_{n+4}^c - y_{n+4}^{p,m}),$$

$$E: f_{n+4} = f(x_{n+4}, y_{n+4}).$$

例 7 将例 6 的初值问题用修正的米尔尼-汉明预测-校正公式计算 y_5 及 y_6 , 初值 y_0, y_1, y_2, y_3 仍用已算出的精确解, 即 $y_0 = 1, y_1 = 1.004\ 837\ 42, y_2 = 1.018\ 730\ 75, y_3 = 1.040\ 818\ 22$, 给出计算结果及误差.

解 根据修正的米尔尼-汉明预测-校正公式可得

$$y_5^p = 1.106\ 532\ 99,$$

$$y_5^{p,m} = y_5^p + \frac{112}{121} \times (y_4^c - y_4^p) = 1.106\ 530\ 364.$$

(注: $y_4^p = 1.070\ 322\ 60, y_4^c = 1.070\ 319\ 66$)

$$f(x_5, y_5^{p,m}) = -y_5^{p,m} + x_5 + 1 = 0.393\ 469\ 636,$$

$$\begin{aligned} y_5^c &= \frac{1}{8} \times (9y_4 - y_2) + \frac{0.3}{8} \times [f(x_5, y_5^{p,m}) + 2f_4 - f_3] \\ &= 1.106\ 530\ 419, \end{aligned}$$

$$y_5 = y_5^c - \frac{9}{121} \times (y_5^c - y_5^p) = 1.106\ 530\ 61;$$

$$f_5 = -y_5 + x_5 + 1 = 0.393\ 469\ 39,$$

$$y_6^p = y_2 + \frac{0.4}{3} \times (2f_5 - f_4 + 2f_3) = 1.148\ 813\ 73,$$

$$y_6^{p,m} = y_6^p + \frac{112}{121} \times (y_5^c - y_5^p) = 1.148\ 811\ 35,$$

$$f_6^{p,m} = -y_6^{p,m} + x_6 + 1 = 0.451\ 188\ 65,$$

$$y_6^c = \frac{1}{8} \times (9y_5 - y_3) + \frac{0.3}{8} \times (f_6^{p,m} + 2f_5 - f_4) = 1.148\ 811\ 44,$$

$$y_6^c = y_5^c - \frac{9}{121} \times (y_5^c - y_5^p) = 1.14881161.$$

误差

$$|y(x_5) - y_5| = 4.97 \times 10^{-8}, |y(x_6) - y_6| = 2.61 \times 10^{-8}.$$

从结果看,此方法的误差比四阶阿当姆斯隐式法和四阶汉明方法小,这与理论分析一致.

9.5.6 构造多步法公式的注记和例

前面已指出构造多步法公式有基于数值积分和泰勒展开两种途径,只对能将微分方程(1.1)转化为等价的积分方程的情形方可利用数值积分方法建立多步法公式,它是有局限性的.即前种途径只对部分方法适用.而用泰勒展开则可构造任意多步法公式,其做法是根据多步法公式的形式,直接在 x_n 处做泰勒展开即可.不必套用系数公式(5.4)确定多步法(5.1)的系数 a_i 及 b_i ($i=0, 1, \dots, k$),因为多步法公式不一定如(5.1)的形式.另外,套用公式容易记错.具体做法见下面例子.

例 8 解初值问题 $y = f(x, y)$, $y(x_0) = y_0$. 用显式二步法 $y_{n+1} = y_n + a_0 y_{n-1} + h(f_n + b_1 f_{n-1})$, 其中 $f_n = f(x_n, y_n)$, $f_{n-1} = f(x_{n-1}, y_{n-1})$. 试确定参数 a_0, b_1, a_0, b_1 使方法阶数尽可能高,并求局部截断误差.

解 本题仍根据局部截断误差定义,用泰勒展开确定参数满足的方程.由于

$$\begin{aligned} T_{n+1} &= y(x_n + h) - y(x_n) - a_0 y(x_n - h) \\ &\quad - h[f(x_n) + b_1 f(x_n - h)] \\ &= y(x_n) + hy'(x_n) + \frac{h^2}{2}y''(x_n) + \frac{h^3}{3!}y'''(x_n) \\ &\quad + \frac{h^4}{4!}y^{(4)}(x_n) + O(h^5) - y(x_n) \\ &\quad - a_0 y(x_n) - hy(x_n) + \frac{h^2}{2}y(x_n) \end{aligned}$$

$$\begin{aligned}
& - \frac{h^3}{3!} y''(x_n) + \frac{h^4}{4!} y^{(4)}(x_n) + O(h^5) \\
& - \alpha_0 h y'(x_n) - \alpha_1 h y(x_n) - h y(x_n) \\
& + \frac{h^2}{2} y''(x_n) - \frac{h^3}{3!} y^{(4)}(x_n) + O(h^4) \\
= & (1 - \alpha_0 - \alpha_1) y(x_n) + (1 + \alpha_1 - \alpha_0 - \alpha_1) h y'(x_n) \\
+ & \frac{1}{2} - \frac{1}{2} \alpha_1 + \alpha_1 h^2 y(x_n) + \frac{1}{6} + \frac{1}{6} \alpha_1 - \frac{1}{2} \alpha_1 h^3 \\
& \cdot y''(x_n) + \frac{1}{24} - \frac{1}{24} \alpha_1 + \frac{1}{6} \alpha_1 h^4 y^{(4)}(x_n) + O(h^5)
\end{aligned}$$

为求参数 $\alpha_0, \alpha_1, \alpha_0, \alpha_1$ 使方法阶数尽量高, 可令

$$\begin{aligned}
1 - \alpha_0 - \alpha_1 &= 0, \quad 1 + \alpha_1 - \alpha_0 - \alpha_1 = 0, \\
\frac{1}{2} - \frac{1}{2} \alpha_1 + \alpha_1 &= 0, \quad \frac{1}{6} + \frac{1}{6} \alpha_1 - \frac{1}{2} \alpha_1 = 0.
\end{aligned}$$

即得方程组

$$\begin{aligned}
\alpha_0 + \alpha_1 &= 1, \\
-\alpha_1 + \alpha_0 + \alpha_1 &= 1, \\
\alpha_1 - 2\alpha_1 &= 1, \\
-\alpha_1 + 3\alpha_1 &= 1.
\end{aligned}$$

解得 $\alpha_0 = -4, \alpha_1 = 5, \alpha_0 = 4, \alpha_1 = 2$, 此时公式为三阶, 而且

$$T_{n+1} = \frac{1}{6} h^4 y^{(4)}(x_n) + O(h^5)$$

即为所求局部截断误差. 而所得二步法为

$$y_{n+1} = -4y_n + 5y_{n-1} + 2h(2f_n + f_{n-1}).$$

例 9 证明存在 λ 的一个值, 使线性多步法

$$y_{n+1} + (\lambda y_n - y_{n-1}) - y_{n-2} = \frac{1}{2}(3 + \lambda)h(f_n + f_{n-1})$$

是四阶的.

证明 只要证明局部截断误差 $T_{n+1} = O(h^5)$, 则方法为四阶. 仍用泰勒展开, 由于

$$\begin{aligned}
T_{n+1} &= y(x_n + h) + [y(x_n) - y(x_n - h)] + y(x_n - 2h) \\
&\quad - \frac{1}{2}(3 +)h[y(x_n) + y(x_n - h)] \\
&= y(x_n) + hy(x_n) + \frac{h^2}{2}y(x_n) + \frac{h^3}{3!}y(x_n) + \frac{h^4}{4!}y^{(4)}(x_n) \\
&\quad + O(h^5) - (-h)y(x_n) + \frac{h^2}{2}y(x_n) - \frac{h^3}{3!}y(x_n) \\
&\quad + \frac{h^4}{4!}y^{(4)}(x_n) + O(h^5) - y(x_n) - 2hy(x_n) \\
&\quad + \frac{(2h)^2}{2}y(x_n) - \frac{(2h)^3}{3!}y(x_n) + \frac{(2h)^4}{4!}y^{(4)}(x_n) + O(h^5) \\
&\quad - \frac{h}{2}(3 +)y(x_n) + y(x_n) - hy(x_n) + \frac{h^2}{2}y(x_n) \\
&\quad - \frac{h^3}{3!}y^{(4)}(x_n) + O(h^4) \\
&= [1 + + 2 - (3 +)]hy(x_n) + \frac{1}{2} - \frac{1}{2} - 2 \\
&\quad + \frac{1}{2}(3 +)h^2y(x_n) + \frac{1}{6} + \frac{1}{6} + \frac{4}{3} \\
&\quad - \frac{1}{4}(3 +)h^3y(x_n) + \frac{1}{24} - \frac{1}{24} - \frac{2}{3} + \frac{1}{12}(3 +) \\
&\quad \times h^4y^{(4)}(x_n) + O(h^5) \\
&= \frac{3}{4} - \frac{1}{12}h^3y(x_n) + \frac{1}{24}(-9 +)h^4y^{(4)}(x_n) + O(h^5).
\end{aligned}$$

当 = 9 时, $T_{n+1} = O(h^5)$, 故方法是四阶的.

9.6 方程组和高阶方程

9.6.1 一阶方程组

前面我们研究了单个方程 $y = f$ 的数值解法, 只要把 y 和 f

理解为向量,那么,所提供的各种计算公式即可应用到一阶方程组的情形.

考察一阶方程组

$$y_i = f_i(x, y_1, y_2, \dots, y_N) \quad (i = 1, 2, \dots, N)$$

的初值问题,初始条件给为

$$y_i(x_0) = y_i^0 \quad (i = 1, 2, \dots, N).$$

若采用向量的记号,记

$$\mathbf{y} = (y_1, y_2, \dots, y_N)^T,$$

$$\dot{\mathbf{y}} = (y_1^0, y_2^0, \dots, y_N^0)^T,$$

$$\mathbf{f} = (f_1, f_2, \dots, f_N)^T.$$

则上述方程组的初值问题可表示为

$$\begin{aligned} \mathbf{y} &= \mathbf{f}(x, \mathbf{y}), \\ \mathbf{y}(x_0) &= \mathbf{y}_0. \end{aligned} \tag{6.1}$$

求解这一初值问题的四阶龙格-库塔公式为

$$\mathbf{y}_{n+1} = \mathbf{y}_n + \frac{h}{6}(\mathbf{k}_1 + 2\mathbf{k}_2 + 2\mathbf{k}_3 + \mathbf{k}_4),$$

式中

$$\mathbf{k}_1 = \mathbf{f}(x_n, \mathbf{y}_n),$$

$$\mathbf{k}_2 = \mathbf{f}(x_n + \frac{h}{2}, \mathbf{y}_n + \frac{h}{2}\mathbf{k}_1),$$

$$\mathbf{k}_3 = \mathbf{f}(x_n + \frac{h}{2}, \mathbf{y}_n + \frac{h}{2}\mathbf{k}_2),$$

$$\mathbf{k}_4 = \mathbf{f}(x_n + h, \mathbf{y}_n + h\mathbf{k}_3).$$

或表示为

$$\begin{aligned} y_{i,n+1} &= y_{i,n} + \frac{h}{6}(K_1 + 2K_2 + 2K_3 + K_4) \\ &\quad (i = 1, 2, \dots, N), \end{aligned}$$

其中

$$K_1 = f_i(x_n, y_{1,n}, y_{2,n}, \dots, y_{N,n}),$$

$$\begin{aligned}
 K_2 &= f_i(x_n + \frac{h}{2}, y_{1n} + \frac{h}{2}K_{11}, y_{2n} + \frac{h}{2}K_{21}, \dots, y_{Nn} + \frac{h}{2}K_{N1}), \\
 K_3 &= f_i(x_n + \frac{h}{2}, y_{1n} + \frac{h}{2}K_{12}, y_{2n} + \frac{h}{2}K_{22}, \dots, y_{Nn} + \frac{h}{2}K_{N2}), \\
 K_4 &= f_i(x_n + h, y_{1n} + hK_{13}, y_{2n} + hK_{23}, \dots, y_{Nn} + hK_{N3}).
 \end{aligned}$$

这里 y_{in} 是第 i 个因变量 $y_i(x)$ 在节点 $x_n = x_0 + nh$ 的近似值.

为了帮助理解这一公式的计算过程, 我们再考察两个方程的特殊情形:

$$\begin{aligned}
 y &= f(x, y, z), \\
 z &= g(x, y, z), \\
 y(x_0) &= y_0, \\
 z(x_0) &= z_0.
 \end{aligned}$$

这时四阶龙格 - 库塔公式具有形式

$$\begin{aligned}
 y_{n+1} &= y_n + \frac{h}{6}(K_1 + 2K_2 + 2K_3 + K_4), \\
 z_{n+1} &= z_n + \frac{h}{6}(L_1 + 2L_2 + 2L_3 + L_4).
 \end{aligned} \tag{6.2}$$

其中

$$\begin{aligned}
 K_1 &= f(x_n, y_n, z_n), \\
 K_2 &= f(x_n + \frac{h}{2}, y_n + \frac{h}{2}K_1, z_n + \frac{h}{2}L_1), \\
 K_3 &= f(x_n + \frac{h}{2}, y_n + \frac{h}{2}K_2, z_n + \frac{h}{2}L_2), \\
 K_4 &= f(x_n + h, y_n + hK_3, z_n + hL_3), \\
 L_1 &= g(x_n, y_n, z_n), \\
 L_2 &= g(x_n + \frac{h}{2}, y_n + \frac{h}{2}K_1, z_n + \frac{h}{2}L_1), \\
 L_3 &= g(x_n + \frac{h}{2}, y_n + \frac{h}{2}K_2, z_n + \frac{h}{2}L_2), \\
 L_4 &= g(x_n + h, y_n + hK_3, z_n + hL_3).
 \end{aligned} \tag{6.3}$$

这是一步法, 利用节点 x_n 上的值 y_n, z_n , 由(6.3)式顺序计算 $K_1, L_1, K_2, L_2, K_3, L_3, K_4, L_4$, 然后代入(6.2)式即可求得节点 x_{n+1} 上的 y_{n+1}, z_{n+1} .

9.6.2 化高阶方程为一阶方程组

关于高阶微分方程(或方程组)的初值问题, 原则上总可以归结为一阶方程组来求解. 例如, 考察下列 m 阶微分方程

$$y^{(m)} = f(x, y, y', \dots, y^{(m-1)}), \quad (6.4)$$

初始条件为

$$y(x_0) = y_0, y'(x_0) = y_1, \dots, y^{(m-1)}(x_0) = y_0^{(m-1)}. \quad (6.5)$$

只要引进新的变量

$$y_1 = y, y_2 = y', \dots, y_m = y^{(m-1)},$$

即可将 m 阶方程(6.4)化为如下的一阶方程组:

$$\begin{aligned} y_1 &= y_2, \\ y_2 &= y_3, \\ &\dots\dots \\ y_{m-1} &= y_m, \\ y_m &= f(x, y_1, y_2, \dots, y_m). \end{aligned} \quad (6.6)$$

初始条件(6.5)则相应地化为

$$y_1(x_0) = y_0, y_2(x_0) = y_1, \dots, y_m(x_0) = y_0^{(m-1)}. \quad (6.7)$$

不难证明初值问题(6.4), (6.5)和(6.6), (6.7)是彼此等价的.

特别地, 对于下列二阶方程的初值问题:

$$y = f(x, y, y'),$$

$$y(x_0) = y_0,$$

$$y'(x_0) = y_1.$$

引进新的变量 $z = y$, 即可化为下列一阶方程组的初值问题:

$$\begin{aligned}y &= z, \\z &= f(x, y, z), \\y(x_0) &= y_0, \\z(x_0) &= y_0.\end{aligned}$$

针对这个问题应用四阶龙格-库塔公式(6.2),有

$$\begin{aligned}y_{n+1} &= y_n + \frac{h}{6}(K_1 + 2K_2 + 2K_3 + K_4), \\z_{n+1} &= z_n + \frac{h}{6}(L_1 + 2L_2 + 2L_3 + L_4).\end{aligned}$$

由(6.3)式可得

$$\begin{aligned}K_1 &= z_n, \quad L_1 = f(x_n, y_n, z_n); \\K_2 &= z_n + \frac{h}{2}L_1, \quad L_2 = f(x_n + \frac{h}{2}, y_n + \frac{h}{2}K_1, z_n + \frac{h}{2}L_1); \\K_3 &= z_n + \frac{h}{2}L_2, \quad L_3 = f(x_n + \frac{h}{2}, y_n + \frac{h}{2}K_2, z_n + \frac{h}{2}L_2); \\K_4 &= z_n + hL_3, \quad L_4 = f(x_n + h, y_n + hK_3, z_n + hL_3).\end{aligned}$$

如果消去 K_1, K_2, K_3, K_4 , 则上述格式可表示为

$$\begin{aligned}y_{n+1} &= y_n + hz_n + \frac{h^2}{6}(L_1 + L_2 + L_3), \\z_{n+1} &= z_n + \frac{h}{6}(L_1 + 2L_2 + 2L_3 + L_4).\end{aligned}$$

这里

$$\begin{aligned}L_1 &= f(x_n, y_n, z_n), \\L_2 &= f(x_n + \frac{h}{2}, y_n + \frac{h}{2}z_n, z_n + \frac{h}{2}L_1), \\L_3 &= f(x_n + \frac{h}{2}, y_n + \frac{h}{2}z_n + \frac{h^2}{4}L_1, z_n + \frac{h}{2}L_2), \\L_4 &= f(x_n + h, y_n + hz_n + \frac{h^2}{2}L_2, z_n + hL_3).\end{aligned}$$

9.6.3 刚性方程组

在求解方程组(6.1)时, 经常出现解的分量数量级差别很大的情形, 这给数值求解带来很大困难, 这种问题称为刚性(stiff)问题, 在化学反应、电子网络和自动控制等领域中都是常见的, 先考察以下例子.

给定系统

$$\begin{aligned} u &= -1000.25u + 999.75v + 0.5, \\ v &= 999.75u - 1000.25v + 0.5, \\ u(0) &= 1, \\ v(0) &= -1. \end{aligned} \quad (6.8)$$

它可用解析方法求出准确解, 方程右端系数矩阵

$$\mathbf{A} = \begin{pmatrix} -1000.25 & 999.75 \\ 999.75 & -1000.25 \end{pmatrix}$$

的特征值为 $\lambda_1 = -0.5$, $\lambda_2 = -2000$, 方程的准确解为

$$\begin{aligned} u(t) &= -e^{-0.5t} + e^{-2000t} + 1, \\ v(t) &= -e^{-0.5t} - e^{-2000t} + 1. \end{aligned}$$

当 $t = 0$ 时, $u(t) = 1$, $v(t) = -1$ 称为稳态解, u , v 中均含有快变分量 e^{-2000t} 及慢变分量 $e^{-0.5t}$. 对应于 λ_2 的快速衰减的分量在 $t = 0.005$ 秒时已衰减到 $e^{-10} = 0$, 称 $\tau_2 = -\frac{1}{\lambda_2} = \frac{1}{2000} = 0.0005$ 为时间常数.

当 $t = 10\tau_2$ 时快变分量即可被忽略, 而对应于 λ_1 的慢变分量, 它的时间常数 $\tau_1 = -\frac{1}{\lambda_1} = \frac{1}{0.5} = 2$, 它要计算到 $t = 10\tau_1 = 20$ 时, 才能衰减到 $e^{-10} = 0$, 也就是说解 u , v 必须计算到 $t = 20$ 才能达到稳态解. 它表明方程(6.8)的解分量变化速度相差很大, 是一个刚性方程组. 如果用四阶龙格-库塔法求解, 步长选取要满足 $h < -2.78/\tau_2$, 即 $h < -2.78/\tau_2 = 0.00139$, 才能使计算稳定. 而要计算到稳态解至少需要算到 $t = 20$, 则需计算 14 388 步. 这种用小步长计算长

区间的现象是刚性方程数值求解出现的困难, 它是系统本身病态性质引起的.

对一般的线性系统

$$\frac{d\mathbf{y}}{dt} = \mathbf{A}\mathbf{y}(t) + \mathbf{g}(t), \quad (6.9)$$

其中 $\mathbf{y} = (y_1, \dots, y_N)^T \in \mathbf{R}^N$, $\mathbf{g} = (g_1, \dots, g_N)^T \in \mathbf{R}^N$, $\mathbf{A} \in \mathbf{R}^{N \times N}$. 若 \mathbf{A} 的特征值 $\lambda_j = \mu_j + i\nu_j$ ($j = 1, \dots, N$, $i = -1$) 相应的特征向量为 \mathbf{v}_j ($j = 1, \dots, N$), 则方程组(6.9)的通解为

$$\mathbf{y}(t) = \sum_{j=1}^N c_j e^{\mu_j t} \mathbf{v}_j + \mathbf{y}^0(t), \quad (6.10)$$

其中 c_j 为任意常数, 可由初始条件 $y^0(a) = y^0$ 确定, $\mathbf{y}^0(t)$ 为特解.

假定 μ_j 的实部 $\mu_j = \operatorname{Re}(\lambda_j) < 0$, 则当 $t \rightarrow \infty$ 时, $\mathbf{y}(t) \rightarrow \mathbf{y}^0(t)$, $\mathbf{y}^0(t)$ 为稳态解.

定义 8 若线性系统(6.9)中 \mathbf{A} 的特征值 λ_j 满足条件 $\operatorname{Re}(\lambda_j) < 0$ ($j = 1, \dots, N$), 且

$$s = \max_{1 \leq j \leq N} |\operatorname{Re}(\lambda_j)| / \min_{1 \leq j \leq N} |\operatorname{Re}(\lambda_j)| \geq 1,$$

则称系统(6.9)为刚性方程, 称 s 为刚性比.

刚性比 $s \geq 1$ 时, \mathbf{A} 为病态矩阵, 故刚性方程也称病态方程. 通常 $s \geq 10$ 就认为是刚性的. s 越大病态越严重. 方程(6.8)的刚性比 $s=4000$, 故它是刚性的.

对一般非线性方程组(6.1), 可类似定义 8, 将 \mathbf{f} 在点 $(t, \mathbf{y}(t))$ 处线性展开, 记 $\mathbf{J}(t) = \frac{\partial \mathbf{f}}{\partial \mathbf{y}} \in \mathbf{R}^{N \times N}$, 假定 $\mathbf{J}(t)$ 的特征值为 $\lambda_j(t)$, $j = 1, \dots, N$, 于是由定义 8 可知, 当 $\lambda_j(t)$ 满足条件 $\operatorname{Re}(\lambda_j(t)) < 0$ ($j = 1, \dots, N$), 且

$$s(t) = \max_{1 \leq j \leq N} |\operatorname{Re}(\lambda_j(t))| / \min_{1 \leq j \leq N} |\operatorname{Re}(\lambda_j(t))| \geq 1,$$

则称系统(6.1)是刚性的, $s(t)$ 称为方程(6.1)的局部刚性比.

求刚性方程数值解时, 若用步长受限制的方法就将出现小步

长计算大区间的问题,因此最好使用对步长 h 不加限制的方法,如前面已介绍的欧拉后退法及梯形法,即 A-稳定的方法,所谓 A-稳定就是指数值方法的绝对稳定域包含了 $\mu = h$ 平面的左半平面。这种方法当然对步长 h 没有限制,但 A-稳定方法要求太苛刻,Dahlquist 已证明所有显式方法都不是 A-稳定的,而隐式的 A-稳定多步法阶数最高为 2,且以梯形法误差常数为最小。这就表明本章所介绍的方法中能用于解刚性方程的方法很少。通常求解刚性方程的高阶线性多步法是吉尔(Gear)方法,还有隐式龙格-库塔法(见文献[21]),这些方法都有现成的数学软件可供使用。本书不再介绍。

评注

本章研究求解常微分方程初值问题的数值方法。构造数值方法主要有两条途径:基于数值积分的构造方法和基于泰勒展开的构造方法。后一种方法更灵活,也更具有一般性。泰勒展开方法还有一个优点,它在构造差分公式的同时可以得到关于截断误差的估计。

基于泰勒展开构造出的四阶龙格-库塔方法(见 3.3 节)则是计算机上的常用算法。四阶龙格-库塔方法的优点是精度高,程序简单,计算过程稳定,并且易于调节步长。

四阶龙格-库塔方法也有不足之处:它要求函数 $f(x, y)$ 具有较高的光滑性。如果 $f(x, y)$ 的光滑性差,那么,它的精度可能还不如欧拉公式(2.1 节)或改进的欧拉公式(2.4 节)。

四阶龙格-库塔方法的另一个缺点是计算量比较大,需要耗费较多的机器时间(每一步需四次计算函数 $f(x, y)$ 的值)。相比之下,汉明方法(5.4 节)可以节省计算量(每一步只需两次计算函数 $f(x, y)$ 的值)。但汉明方法是一种四步法,它不是自开始的,需要

借助于四阶龙格-库塔方法提供开始值.

对数值方法的分析还涉及到局部截断误差, 整体误差, 收敛性, 相容性及稳定性等概念, 特别是有关绝对稳定性的讨论, 它涉及计算时步长的选取, 如步长选得不合适, 舍入误差恶性增长, 结果完全错误. 本章只对单步法的收敛性和稳定性做了讨论. 对多步法的有关理论没有涉及, 读者可参阅文献[2]及[21].

刚性方程组是具有重要应用价值的问题, 它的理论和解法很多, 并且还在进一步发展, 详细介绍可参阅文献[21]和[22].

习 题

1. 用欧拉法解初值问题

$$y = x^2 + 100y^2, y(0) = 0.$$

取步长 $h=0.1$, 计算到 $x=0.3$ (保留到小数点后 4 位).

2. 用改进欧拉法和梯形法解初值问题

$$y = x^2 + x - y, y(0) = 0.$$

取步长 $h=0.1$, 计算到 $x=0.5$, 并与准确解 $y = -e^{-x} + x^2 - x + 1$ 相比较.

3. 用梯形方法解初值问题

$$y + y = 0,$$

$$y(0) = 1.$$

证明其近似解为

$$y_h = \frac{2-h}{2+h}^n,$$

并证明当 $h \rightarrow 0$ 时, 它收敛于原初值问题的准确解 $y = e^{-x}$.

4. 利用欧拉方法计算积分

$$\int_0^x e^t dt$$

在点 $x=0.5, 1, 1.5, 2$ 的近似值.

5. 取 $h=0.2$, 用四阶经典的龙格-库塔方法求解下列初值问题:

$$1) \quad y = x + y, \quad 0 < x < 1; \\ y(0) = 1;$$

2) $y' = 3y'(1+x), \quad 0 < x < 1;$
 $y(0) = 1.$

6. 证明对任意参数 t , 下列龙格-库塔公式是二阶的:

$$\begin{aligned} y_{n+1} &= y_n + \frac{h}{2}(K_2 + K_3), \\ K_1 &= f(x_n, y_n), \\ K_2 &= f(x_n + th, y_n + thK_1), \\ K_3 &= f(x_n + (1-t)h, y_n + (1-t)hK_1). \end{aligned}$$

7. 证明中点公式(3.10)是二阶的, 并求其绝对稳定区间.

8. 对于初值问题

$$y' = -100(y - x^2) + 2x, \quad y(0) = 1.$$

(1) 用欧拉法求解, 步长 h 取什么范围的值, 才能使计算稳定.

(2) 若用四阶龙格-库塔法计算, 步长 h 如何选取?

(3) 若用梯形公式计算, 步长 h 有无限制.

9. 分别用二阶显式阿当姆斯方法和二阶隐式阿当姆斯方法解下列初值问题:

$$y' = 1 - y, \quad y(0) = 0.$$

取 $h=0.2$, $y_0=0$, $y_1=0.181$, 计算 $y(1.0)$ 并与准确解 $y=1-e^{-x}$ 相比较.

10. 证明解 $y = f(x, y)$ 的下列差分公式

$$y_{n+1} = \frac{1}{2}(y_n + y_{n-1}) + \frac{h}{4}(4y_{n+1} - y_n + 3y_{n-1})$$

是二阶的, 并求出截断误差的主项.

11. 试证明线性二步法

$$y_{n+2} + (b-1)y_{n+1} - by_n = \frac{h}{4}[(b+3)f_{n+2} + (3b+1)f_n]$$

当 $b = -1$ 时方法为二阶, 当 $b = -1$ 时方法为三阶.

12. 求方程

$$\begin{aligned} u' &= -10u + 9v, \\ v' &= 10u - 11v. \end{aligned}$$

的刚性比, 用四阶 R-K 方法求解时, 最大步长能取多少?

计算实习题

要求：

(1) 用 Matlab 语言或你熟悉的其他算法语言编程序, 使之尽量具有通用性 .

(2) 上机前充分准备, 复习有关算法, 写出计算步骤, 反复查对程序 .

(3) 完成计算后写出实验报告, 内容包括: 计算机型号及所用算法语言, CPU 时间, 算法步骤叙述, 变量说明, 程序清单, 输出计算结果, 结果分析和小结等 .

(4) 根据教师要求选做下列计算实习题中的 3 或 5 个题目 .

1 . 用三次样条插值的三弯矩法, 编制第一与第二种边界条件的程序 . 设已知数据如下:

x_i	0 . 2	0 . 4	0 . 6	0 . 8	1 . 0
$f(x_i)$	0 . 9798652	0 . 9177710	0 . 8080348	0 . 6386093	0 . 3843735

求 $f(x)$ 的三次样条插值函数 $S(x)$, 满足:

(1) 自然边界条件 $S(0 . 2) = S(1 . 0) = 0$;

(2) 第一种边界条件 $S(0 . 2) = 0 . 20271$, $S(1 . 0) = 1 . 55741$.
要求输出用追赶法解出的弯矩向量 (M_0, M_1, \dots, M_4) 和 $S(0 . 2 + 0 . 1 i) (i = 0, 1, \dots, 8)$ 的值 . 并画出 $y = S(x)$ 的图形 .

2 . 编制以离散点 $\{ x_i \}_{i=0}^m$ 的正交多项式 $\{ P_k(x) \}$ 为基的最小二乘拟合程序, 并用于对下列数据做三次多项式最小二乘拟合 .

x_i	- 1 .0	- 0 .5	0 .0	0 .5	1 .0	1 .5	2 .0
y_i	- 4 .447	- 0 .452	0 .551	0 .048	- 0 .447	0 .549	4 .552

取权 (x_i) 1, 求出拟合曲线 $y = S(x) = \sum_{k=0}^3 P_k(x) = a_0 + a_1 x + a_2 x^2 + a_3 x^3$, 输出 $P_k(x)$, a_k ($k = 0, 1, 2, 3$) 及平方误差 $\sum_{i=1}^n (y_i - S(x_i))^2$, 并画出 $y = S(x)$ 的图形.

3. 给出积分

$$(i) I = \int_0^2 x^2 e^{-x^2} dx, \quad (ii) I = \int_0^{\frac{3}{4}} \cot x dx,$$

$$(iii) I = \int_2^3 \frac{1}{x^2 - 1} dx.$$

实验要求:

(1) 用龙贝格求积计算上述积分 I 的值, 要求到 $|T_{k+1}^{(k+1)} - T_k^{(k)}| < 10^{-6}$ 时结束, 输出 T 表及 I 的近似值.

(2) 用 5 点高斯求积公式及复化 3 点高斯求积公式计算上述积分, 并输出 I 的近似值.

(3) 分析比较各种计算结果.

4. 比较求一阶导数的数值方法, 给出函数

$$f(x) = e^x, \quad x \in [0.5, 2].$$

实验要求: 利用某距离点函数值, 必要时给定端点导数值, 分别用中心差分, 数值积分求导, 三次样条求导和理查森外推计算 $f(x)$ 的一阶导数, 分析, 比较各种方法的效果, 说明精度与步长 h 的关系.

5. 给定方程组

$$(i) \begin{array}{ccccc} 3 .01 & 6 .03 & 1 .99 & x_1 & 1 \\ 1 .27 & 4 .16 & -1 .23 & x_2 & = 1 ; \\ 0 .987 & -4 .81 & 9 .34 & x_3 & 1 \end{array}$$

$$(ii) \begin{array}{ccccccccc} 10 & -7 & & 0 & 1 & x_1 & 8 \\ -3 & 2 & 0.99999 & 6 & 2 & x_2 & = & 5.900001 \\ 5 & -1 & & 5 & -1 & x_3 & & 5 \\ 2 & 1 & & 0 & 2 & x_4 & & 1 \end{array} .$$

实验要求:

(1) 用 LU 分解和列主元高斯消去法求解上述两个方程组。输出 $\mathbf{Ax} = \mathbf{b}$ 中矩阵 \mathbf{A} 及向量 \mathbf{b} , $\mathbf{A} = LU$ 分解的 \mathbf{L} 与 \mathbf{U} , $\det \mathbf{A}$ 及解向量 \mathbf{x} 。

(2) 将方程组(i)中系数 3.01 改为 3.00, 0.987 改为 0.990。用列主元高斯消去法求解, 输出列主元行交换次序、解向量 \mathbf{x} 及 $\det \mathbf{A}$, 并与(1)中结果比较。

(3) 将方程组(ii)中的 2.099999 改为 2.1, 5.900001 改为 5.9。用列主元高斯消去法求解, 输出解向量 \mathbf{x} 及 $\det \mathbf{A}$, 并与(1)中结果比较。

6. 研究解线性方程组 $\mathbf{Ax} = \mathbf{b}$ 迭代法收敛速度。给定 \mathbf{A}
 $\mathbf{R}^{20 \times 20}$ 为五对角矩阵

$$\mathbf{A} = \begin{array}{ccccccccc} 3 & -1/2 & & -1/4 & & & & & \\ -1/2 & 3 & -1/2 & & -1/4 & & & & \\ -1/4 & -1/2 & 3 & -1/2 & & -1/4 & & & \\ & & w & w & w & w & & & \\ & -1/4 & -1/2 & 3 & -1/2 & -1/4 & & & \\ & -1/4 & -1/2 & 3 & -1/2 & -1/4 & & & \\ & & -1/4 & -1/2 & 3 & & & & \end{array} .$$

实验要求:

(1) 选取不同的初始向量 $\mathbf{x}^{(0)}$ 及右端项向量 \mathbf{b} , 给定迭代误差要求, 用雅可比迭代和高斯-塞德尔迭代法求解, 观察得到的序列是否收敛? 若收敛, 记录迭代次数, 分析计算结果并得出你的结论。

(2) 用 SOR 迭代法求上述方程组的解, 松弛系数 取 $1 < \omega < 2$ 的不同值, 在 $\|\mathbf{x}^{(k)} - \mathbf{x}^{(k+1)}\| < 10^{-5}$ 时停止迭代. 记录迭代次数. 分析计算结果并得出你的结论.

7. 求非线性方程及方程组的根, 准确到 10^{-6} . 给定方程分别为

$$(i) \quad x^2 - 3x + 2 - e^x = 0.$$

$$(ii) \quad \begin{aligned} 3x_1^2 - x_2^2 &= 0, \\ 3x_1 x_2^2 - x_1^3 - 1 &= 0, \end{aligned} \quad \mathbf{x}^{(0)} = (1, 1)^T.$$

实验要求:

(1) 用你自己设计的一种线性收敛迭代法求方程(i)的根, 然后再用斯蒂芬森加速迭代计算.

(2) 用牛顿法求方程(i)的根, 输出迭代初值, 各次迭代值及迭代次数, 并与(1)的结果比较.

(3) 用牛顿法求(ii)的解, 输出迭代次数及解向量 \mathbf{x} 的近似.

8. 用 QR 算法求矩阵特征值:

$$(i) \quad (\mathbf{A}) = \begin{matrix} & 2 & 3 & 4 & 5 & 6 \\ 6 & 2 & 1 & & & \\ 2 & 3 & 1 & & & \\ 1 & 1 & 1 & & & \end{matrix}, \quad (ii) \quad (\mathbf{H}) = \begin{matrix} & 4 & 4 & 5 & 6 & 7 \\ 0 & 3 & 6 & 7 & 8 & . \\ 0 & 0 & 2 & 8 & 9 & \\ 0 & 0 & 0 & 1 & 0 & \end{matrix}.$$

实验要求:

(1) 根据 QR 算法原理编制求(i)及(ii)中矩阵全部特征值的程序并输出计算结果(要求误差 $< 10^{-5}$).

(2) 直接用现有数学软件求(i), (ii)的全部特征值, 并与(1)的结果比较.

9. 求初值问题的数值解, 给定初值问题为

$$(i) \quad y' = \frac{1}{x^2} - \frac{y}{x}, \quad 1 \leq x \leq 2,$$

$$y(1) = 1.$$

$$y' = -50y + 50x^2 + 2x, \quad 0 \leq x \leq 1,$$

(ii) $y(0) = \frac{1}{3}$.

实验要求:

(1) 用改进欧拉法(取 $h = 0.05$)及四阶 R-K 方法(取 $h = 0.1$)求(i)的数值解, 并输出 $x_i = 0.1i$ ($i = 0, 1, \dots, 10$) 的数值解 y_i .

(2) 用经典四阶 R-K 方法解(ii), 步长 h 分别取为 $h = 0.1, 0.025, 0.01$ 计算, 并打印 $x_i = 0.1i$ ($i = 0, 1, \dots, 10$) 各点的数值解 y_i 及准确解 $y(x)$, 并分析结果.

初值问题(ii)的准确解 $y(x) = \frac{1}{3}e^{-50x} + x^2$

附录 并行算法及其基本概念

由于并行计算机的逐步普及,很多高校都配备了大型的多处理机,学习并掌握并行算法对进行大规模的科学与工程计算是十分必要和有益的.为使读者对并行算法有初步了解,特在附录中对并行算法的基本概念做简单介绍,需要进一步了解学习的可参看文献[23]及[24].

1 并行算法及其分类

传统计算机采用冯·诺伊曼(Von Neumann)结构,其特点是每一时刻按一条指令加工处理一个或一对数据,这类计算机称作单指令流单数据流系统,简称SISD(Single Instruction Stream Single Data Stream)型,只有一个进程的算法,就是传统的串行算法,适用于这类计算机.并行计算机与此不同,它每一时刻可按多条指令处理多个数据,面向并行计算机具有2个以上进程的算法称为并行算法.

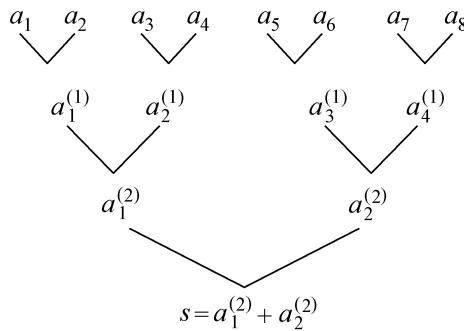
例1 求 N 个数 a_1, a_2, \dots, a_n 的和

$$s = \sum_{i=1}^N a_i. \quad (1)$$

一个进程的串行算法是累加算法

$$s_1 = a_1, \quad s_k = s_{k-1} + a_k, \quad k = 2, \dots, N.$$

$s_N = s$ 即为所求,在累加过程中所给和数规模逐次减1,显然它不适合于并行计算.附图1给出 $N=8$ 时的一种并行算法,称为扇入加法.它可组成4个进程 P_1, P_2, P_3, P_4 如下:



附图 1 扇入加法示意图

P_1 $a_1^{(1)} = a_1 + a_2$ $a_1^{(2)} = a_1^{(1)} + a_2^{(1)}$ $s = a_1^{(2)} + a_2^{(2)}$	P_2 $a_2^{(1)} = a_3 + a_4$ $a_2^{(2)} = a_3^{(1)} + a_4^{(1)}$	P_3 $a_3^{(1)} = a_5 + a_6$ $a_3^{(2)} = a_5^{(1)} + a_6^{(1)}$	P_4 $a_4^{(1)} = a_7 + a_8$
--	---	---	----------------------------------

在有 4 台处理机的并行系统中用 $3 = 1b8$ 级完成 8 个数求和 . 第 1 级并行做 4 个加法 , 第 2 级并行做 2 个加法 , 第 3 级做 1 个加法 .

并行机必须是包含 2 台以上处理机 , 它有不同类型 . 用一个控制器控制多台处理机 , 在每一时刻都执行同一指令对单个或一对数组进行同样加工 , 这类并行机称为单指令流多数据流系统 , 简称 SIMD(Single Instruction Stream Multiple Data Stream) 机 . 另一类并行机由多个控制器分别控制多台处理机 , 各处理机在自己的指令控制下处理自己的数据流 , 称为多指令流多数据流系统 , 简称 MIMD(Multiple Instruction Stream Multiple Data Stream) 机 .

按指令流单个还是多个并行算法可分为两大类 : **SIMD** 算法与 **MIMD** 算法 .

SIMD 算法的基本特征是 :

(1) k 个进程由 1 个指令流控制 , 一个并行步仅由一条指令控制 , 故每个并行步所含的操作必须完全相同 .

(2) 允许并行步中操作的个数在 2 至 k 之间有任意性 , 且有含 k 个操作的并行步 .

(3) 能在一个具有 k 台处理机的 SIMD 系统中实现 .

例 1 给出的扇入加法就是 SIMD 算法 .

MIMD 算法的基本特征是:

(1) k 个进程由 k 个指令流控制 . 故每个并行步所含各操作可两两相异, 而且存在包含不同操作的并行步 .

(2) 同 SIMD 算法特征(2) .

(3) 能在一个具有 k 台处理机的 MIMD 系统中实现 .

按照进程之间是否需要同步也可将并行算法分为同步算法与异步算法 .

同步算法 是指在 k 个进程的算法中有些进程的若干算法必须在另一些进程的某些算法之后执行, 为此有些进程可能出现在计算之前或计算之间的等待阶段, 同步算法可在具有 k 台处理机的 SIMD 系统或 MIMD 系统中实现 .

异步算法 k 个进程间有信息联系但不须同步, 它只能在一个具有 k 台处理机的 MIMD 系统中实现 . 因此异步算法一定是 MIMD 算法 使用这种算法时各次执行的实际过程可能互不重复 . 同步算法既可以是 SIMD 算法也可以是 MIMD 算法, 而 SIMD 算法一定是同步算法 . 进程间的同步通过单指令流的控制来保证 . 这样并行算法可分成 3 类, 即 SIMD 算法, 同步 MIMD 算法和异步算法 .

例 2 对 $f(x) = 0$ 求根的牛顿法(第 7 章(4.2))分别设计两个进程的同步及异步算法 .

解 同步算法: 可将每次迭代分成三个计算元: 分别计算

$f(x_k)$, f'_k , $f(x_k) - \frac{f_k}{f'_k}$ 及检验精度 . 将计算分为

两个进程 P_1 及 P_2 , 假定计算 $f(x_i)$ 比 $f'(x_i)$ 更花时间, 则两个进程或其中一个必出现等待继续计算所需数据的情况, 附图 2 给出两种同步 MIMD 算法示意图 .

异步算法: 可引进 3 个公用变量 t_1 , t_2 , t_3 分别表示 $f(x)$,

$f(x)$ 及 x 在计算中的当前值, 仍假定计算 $f(x)$ 比计算 $f(x)$ 更花时间, 附图 3 给出两个进程的异步算法, 进程 P_1 更新 t_1 与 t_3 , 且检验根的近似是否满足精度要求, 进程 P_2 只更新 t_2 , 假定变量初值 $t_1 = 0$, $t_2 = c - 0$, $t_3 = x_0$ 前三个近似值

$$\begin{aligned}x_1 &= x_0 - f(x_0) / f'(x_0), \\x_2 &= x_1 - f(x_1) / f'(x_1), \\x_3 &= x_2 - f(x_2) / f'(x_2).\end{aligned}$$

P_1	P_2	P_1	P_2
$f(x_0) \quad f_0$	$f(x_0) \quad f_0$	$f(x_0) \quad f_0$	$f(x_0) \quad f_0$
等待 f_0		等待 x_1	
$x_0 - f_0 / f_0 \quad x_1$ 检验精度	等待 x_1		$x_0 - f_0 / f_0 \quad x_1$ 检验精度
$f(x_1) \quad f_1$	$f(x_1) \quad f_1$	$f(x_1) \quad f_1$	$f(x_1) \quad f_1$
等待 f_1		等待 x_2	
$x_1 - f_1 / f_1 \quad x_2$ 检验精度	等待 x_2		$x_1 - f_1 / f_1 \quad x_2$ 检验精度

附图 2 MIMD 算法示意图

P_1	P_2
$f(t_3) \quad t_1$	
$t_3 - t_1 / t_2 \quad t_3$, 检验	$f(t_3) \quad t_2$
$f(t_3) \quad t_1$	
$t_3 - t_1 / t_2 \quad t_3$, 检验	$f(t_3) \quad t_2$
$f(t_3) \quad t_1$	
$t_3 - t_1 / t_2 \quad t_3$, 检验	$f(t_3) \quad t_2$

附图 3 异步算法示意图

其一般关系如下:

$$x_{k+1} = x_k - f(x_k) / f'(x_k), \quad k = 0, 1, \dots, \quad j < k. \quad (2)$$

当 $k = j$ 时, 这与传统牛顿法(4.2)不同, 它是一种混乱迭代, 其收敛性要另外证明, 由于“计算 $f(x)$ 比 $f(x_j)$ 更花时间”, 所以, 只要 $f(x_j)$ 计算出新值就用于迭代, 对异步算法, 由于不用等待故它更节省机时.

2 并行算法基本概念

评价和分析并行算法, 主要应关注它的计算复杂性, 即算法的运行时间(时间复杂性)与所提供的处理机台数(空间复杂性), 并行处理的基本思想是用增加处理机台数来换取运算时间的节省, 为了评价并行算法我们先引进一些基本概念.

定义 1 一个算法的并行度是指算法中能用一个运算步(并行)完成的运算个数. 假设算法运算个数为 r , 利用 s 个运算步完成, 则 r/s 称为平均并行度.

如例 1 的 N 个数求和(1), 若用串行算法需要 $N - 1$ 个加法步, 每步一个加法, 并行度为 1, 平均并行度 $(N - 1)/(N - 1)$ 也等于 1, 如果用扇入加法, 对 $N = 2^n$ 个数求和, 有 $n = \lg N$ 个加法步, 并行性由高到低分布, 逐步减半, 运算个数 $r = \frac{N}{2} + \frac{N}{2^2} + \dots + 1 = 2^n - 1 = N - 1$, 运算步 $s = n = \lg N$, 其平均并行度为

$$\frac{r}{s} = \frac{N - 1}{\lg N} = O\left(\frac{N}{\lg N}\right). \quad (3)$$

N 个数求和的 $N - 1$ 个加法不能用一个并行步完成, 即不能以 $N - 1$ 为并行度, 但用扇入加法可把并行度降为 $\frac{N - 1}{\lg N}$.

并行度和平均并行度是算法内在并行性的一种度量, 不依赖于并行系统中处理机台数, 当然处理机台数会影响完成计算所需时间.

处理机台数充分多时的最少运行时间称作算法的时间界, 而

算法的运行时间达到时间界时需要提供的处理机台数称作处理机台数界, 上述 N 个数求和的时间界为 $T = \lg N$, 而处理机台数界为 $P = \frac{N}{2}$.

下面考察由 P 台处理机组成的并行系统, 假定每台处理机的算术运算的操作时间相同, 我们引进加速比和效率的概念作为并行化的又一种度量.

定义 2 设 T_1 为串行算法在单处理机上运行时间, T_P 为并行算法在 P 台处理机的系统上运行时间, 则一个并行算法的加速比定义为

$$S_P = \frac{T_1}{T_P}, \quad (4)$$

该算法的效率定义为

$$E_P = \frac{S_P}{P}. \quad (5)$$

加速比和效率是评价并行算法的重要指标, 研究并行算法的目标是达到尽可能大的加速比, 理想上 $S_P = P$, 这时并行效率达到 $E_P = 1$, 但一般不可能达到理想的加速比 $S_P = P$. 这是因为:

- (1) 算法缺乏必要的并行度,
- (2) 在并行系统上没有达到完全的负荷均衡,
- (3) 通讯、存储争用及同步时间延迟等.

3 并行算法设计与二分技术

“分而治之”是并行算法设计的重要原则, 其基本思想是把问题依次划分为可以独立完成的较小问题, 将规模逐次减半的二分技术是并行算法设计的一种基本技术, 例 1 给出求和的扇入加法就是一种二分技术, 设 $N = 2^n$ 将所给和式下标为奇偶的对应项两两合并, 得

$$s = \sum_{i=1}^{N/2} (a_{i-1} + a_i).$$

这样,若令

$$a_i^{(1)} = a_{i-1} + a_i, \quad i = 1, 2, \dots, N/2,$$

则所给和式的规模被压缩了一半

$$s = \sum_{i=1}^{N/2} a_i^{(1)}.$$

重复施行这种规模减半的二分手续,二分 k 次后和式的项数被缩减成 $N/2^k$

$$s = \sum_{i=1}^{N/2^k} a_i^{(k)}.$$

式中

$$a_i^{(k)} = a_{i-1}^{(k-1)} + a_i^{(k-1)}, \quad i = 1, 2, \dots, N/2^k. \quad (6)$$

这样二分 $n=\lg N$ 次后, 所给和式最终退化为一项, 从而直接得出所求和值 s . 于是有求和二分算法: 对 $k=1, 2, \dots, n (= \lg N)$ 执行算式(6)得

$$s = a_1^{(n)}.$$

附表 1 就 $N=8$ 的情形具体给出该算法计算流程 .

附表 1 求和二分算法流程

a_1	a_2	a_3	a_4	a_5	a_6	a_7	a_8
$a_1^{(1)}$		$a_2^{(1)}$		$a_3^{(1)}$		$a_4^{(1)}$	
	$a_1^{(2)}$				$a_2^{(2)}$		
				$a_1^{(3)}$			

从以上求和二分算法看到它有以下特点:

- (1) 结构递归, 每步都是同样求和问题 .
- (2) 规模递减, 和式规模逐次减半 .

概括地说,二分算法的设计原理是反复地将所给计算问题加工成规模减半的同类计算问题,直至规模为1时直接得出问题的解.

上述求和的并行算法加速比为 $S = \frac{N}{\ln N}$, 而处理机台数界为

$P = \frac{N}{2}$. 故在具有 $P = \frac{N}{2}$ 台处理机并行系统中其并行效率为

$$E = \frac{2}{\ln N}.$$

下面再讨论多项式求值问题

$$p(x) = \sum_{i=1}^N a_i x^{i-1}. \quad (7)$$

显然 $x=1$ 时就是求和问题(1). 故可仿照求和二分法, 将所给多项式(7)的奇偶项两两合并, 则有

$$\begin{aligned} p = p(x_0) &= \sum_{i=1}^{N/2} (a_{i-1} x_0^{2(i-2)} + a_i x_0^{2(i-1)}) \\ &= \sum_{i=1}^{N/2} (a_{i-1} + a_i x_0) x_0^{2(i-1)}. \end{aligned}$$

这样,若令

$$a^{(1)} = a_{i-1} + a_i x_0, \quad i = 1, 2, \dots, \frac{N}{2}, \quad x_1 = x_0^2,$$

则有

$$p = \sum_{i=1}^{N/2} a_i^{(1)} x_1^{i-1}.$$

这样加工得出的是一个以 x_1 为变元的多项式, 而其规模(项数)是原来的一半, 因此上述手续是一种二分手续. 重复这种手续, 二分 k 次后所给多项式被加工成

$$p = \sum_{i=1}^{N/2^k} d_i^{(k)} x_k^{i-1} \quad (\text{记 } a_i^{(0)} = a_i),$$

这里

$$\begin{aligned} a_i^{(k)} &= a_{i-1}^{(k-1)} + a_i^{(k-1)} x_{k-1}, \quad i = 1, \dots, \frac{N}{2^k}; \\ x_k &= x_{k-1}^2, \quad k = 1, 2, \dots, n. \end{aligned} \tag{8}$$

这样二分 $n = \lg N$ 次, 最终得出的系数 $a^{(n)}$ 即为所求的值 p .

于是得多项式求值的二分算法:

对 $k = 1, 2, \dots, n$ ($= \lg N$) 执行算式(8), 结果有

$$p = a^{(n)}.$$

现在分析算法的效率, 设将(8)式的 $a_i^{(k)}$ 与 x_k 并行计算 (x_k 可表示为 $x_k = 0 + x_{k-1} + x_{k-1}$ 则与 a_k 同步), 则每步 2 次运算 (一乘一加), 上述算法共做 $n = \lg N$ 步, 故其时间界为

$$T = 2 \lg N.$$

而按(8)并行计算第 k 步需处理机 $N/2^k + 1$ 台, 因此处理机台数界

$$P = \max_{1 \leq k \leq n} \frac{N}{2^k} + 1 = \frac{N}{2} + 1.$$

注意到多项式求值的串行算法需做 $T_1 = 2(N - 1)$ 次运算, 易知此算法的加速比与效率为

$$S_p = \frac{N-1}{\lg N} = \frac{N}{\lg N},$$

$$E_p = \frac{2}{\lg N}.$$

上面分析二分法效率是在提供 $\frac{N}{2} + 1$ 台处理机的系统上得到的, 如果处理机只有 $P + 1$ 台, 而 $P < \frac{N}{2}$, 若假定 $N = rP$, $r > 2$ 为正整数, 此时可将所给多项式(7)分成 P 段, 每段含 r 项进行处理.

可利用串行算法来改造或设计并行算法.许多传统的数值计算方法虽然只适合于串行机上计算, 但其中有不少算法包含了可

直接利用的并行性,例如,上述多项式求值的秦九韶算法可用二分技术改造为适合于并行系统的并行算法.用二分技术还可将许多矩阵计算及解线性方程组直接解法改造成适合于并行计算的算法.还有如解线性方程组的雅可比迭代法本身就具有直接的并行性,还有一些迭代法经过重新排序也可改造成适合于并行的算法.另一类就是根据并行算法特点设计具有新思想的新算法,它的出发点仍然是“分而治之”的原理,符合此原理的区域、算子、系统的分裂方法和技术是设计和实现并行处理的重要手段,如解线性与非线性方程的多分裂方法,解微分方程的区域分解法(即 DDM 方法),都属于此类.由于这些方法减小了求解方程的规模,使计算时间减少,因此效率大为提高.此外,异步数值算法基本上是混乱迭代法,是并行算法最富有特色的组成部分之一,混乱迭代不是传统意义的迭代法,在理论上必须做收敛性与收敛速度的分析,混乱迭代与多分裂技术的结合是近十几年才发展起来的新算法,仍有不少问题值得研究.

关于具体数值计算方法的并行算法可参看文献[23].

参 考 文 献

- 1 李庆扬,易大义,王能超 现代数值分析 北京:高等教育出版社,1995
- 2 关治,陆金甫 数值分析基础 北京:高等教育出版社,1998
- 3 冯康等编 数值计算方法 .北京:国防工业出版社,1978
- 4 黄友谦,李岳生 数值逼近 第2版 北京:高等教育出版社 .1987
- 5 徐利治,王仁宏,周蕴时 函数逼近的理论与方法 .上海:上海科学技术出版社,1983
- 6 沈燮昌 多项式最佳逼近的实现 .上海:上海科学技术出版社,1984
- 7 Stoer J . and Bulirsch R . Introduction to Numerical Analysis . New York: Springer-Verlag, 1980
- 8 Wilkinson J H . Rounding Errors in Algebraic Processes . London: H M . Stationery Office, 1963
- 9 Moore R E . Interval Analysis . New Jersey: Prentice-Hall, 1966
- 10 Nürnberger G . Approximation by Spline Functions . Berlin: Springer-Verlag, 1989
- 11 Nussbaumer H J . Fast Fourier Transform and Convolution Algorithms . Berlin: Springer-Verlag , 1981
- 12 赵访熊,李庆扬 富利叶变换滤波在地震勘探数字处理中的应用 .清华大学学报,1978(4)
- 13 Stewart G W 著 .王国荣等译 矩阵计算引论 .上海:上海科学技术出版社,1980
- 14 Drtega J M 著 张丽君等译 数值分析 北京:高等教育出版社,1973
- 15 Varga R S 著 蒋尔雄译 矩阵迭代分析 .上海: 上海科学技术出版社,1966
- 16 Hageman L A, Young D M 著 .蔡大用等译 实用迭代法 .北京:清华大学出版社,1981
- 17 李庆扬,莫孜中,祁力群 非线性方程组的数值解法 北京:科学出版社,1987

-
- 18 清华大学、北京大学计算方法编写组 .计算方法 .北京:科学出版社,1974
 - 19 Golub G H , Van Loan C F 著 .廉庆荣等译 .矩阵计算 .大连:大连理工大学出版社,1988
 - 20 Wilkinson J H 著 .石钟慈等译 .代数特征值问题 .北京:科学出版社,1987
 - 21 李庆扬 .常微分方程数值解法(刚性问题与边值问题) 北京:高等教育出版社,1992
 - 22 袁兆鼎,费景高,刘德贵 .刚性常微分方程初值问题的数值解法 .北京:科学出版社,1987
 - 23 陈景良 .并行算法引论 北京:石油工业出版社,1992
 - 24 李晓梅,蒋增荣 .并行算法 .长沙:湖南科学技术出版社,1992
 - 25 李庆杨,关治,白峰杉 .数值计算原理 .北京:清华大学出版社,2000

部分习题答案

第 1 章

- 1 . ; 2 . $0.02n$; 4 . 1.05×10^{-3} , $0.215, 10^{-5}$; 6 . $\frac{1}{2} \times 10^{-3}$; 9 . 边长误差
 $< 0.005\text{cm}$; 11 . 不稳定, $y_{10}^* = \frac{1}{2} \times 10^8$; 13 . 0.3×10^{-2} , 0.834×10^{-6} .

第 2 章

- 1 . $x^2/6 + 3x/2 - 7/3$; 2 . -0.620219 , -0.616839 ; 3 . 1.06×10^{-8} ;
6 . $h = 0.006$; 7 . $y_n = 2^n$, $y_n = 2^{n-2}$; 14 . 1, 0; 16 . $P(x) = x^2(x-3)^2/4$;
18 . $|R_1(x)| = h^2/4$; 19 . $|R_4(x)| = h^4/16$.

第 3 章

- 1 . $B_1(f, x) = x$, $B_3(f, x) = 1.5x - 0.402x^2 + 0.098x^3$;
4 . (1) $f_0 = 1$, $f_1 = \sqrt{4}$, $f_2 = \sqrt{7}$, (2) $f_0 = \sqrt{2}$,
 $f_1 = \sqrt{4}$, $f_2 = \frac{5}{6}$, (3) $f_m = \frac{1}{2}^{m+n}$, $f_1 =$
 $\frac{n! m!}{(n+m+1)!} f_2 = \frac{(2n)! (2m)!}{[2(n+m)+1]!}$, (4) $f_0 = 4/e$, $f_1 = 5$
 $- 10/e$, $f_2 = 7/4 - 4/e^2$;

6 . (1) 不构成内积, (2) 构成内积;

- 7 . $T_0^*(x) = 1$, $T_1^*(x) = 2x - 1$, $T_2^*(x) = 8x^2 - 8x + 1$, $T_3^*(x) = 32x^3 - 48x^2 + 18x - 1$;
8 . $T_0(x) = 1$, $T_1(x) = x$, $T_2(x) = x^2 - \frac{2}{5}$, $T_3(x) = x^3 - \frac{9}{14}x$;

$$11. P_0(x) = (M+m)/2; \quad 12. 3/4; \quad 13. P_1(x) = \frac{2}{x} + 0.105257;$$

$$14. P_1(x) = (e-1)x + \frac{1}{2}[e - (e-1)\ln(e-1)]; \quad 15. P_3^*(x) = 5x^3 - \frac{5}{4}x^2 + \frac{1}{4}x - \frac{129}{128}; \quad 16. S^*(x) = 0.1171875 + 1.640625x^2 - 0.8203125x^4;$$

$$17. (1) S_l^*(x) = -0.2958x + 1.1410, \quad (2) S_l^*(x) = 0.1878x + 1.6244, \\ (3) S_l^*(x) = -0.24317x + 1.2159, \quad (4) S_l^*(x) = 0.6822x - 0.6371;$$

$$18. S_3^*(x) = 1.5531913x - 0.5622285x^3; \quad 19. S(t) = 22.25376t - 7.855048; \quad 20. y(x) = 0.9726046 + 0.0500351x^2;$$

$$21. y = 5.2151048e^{-\frac{7.4961692}{t}}; \quad 23. R_{22}(x) = 3 - \frac{4}{x+0.5} + \frac{1.25}{x+1.5};$$

$$24. R_{33}(x) = \frac{60x - 7x^3}{60 + 3x^2}; \quad 25. R_{21}(x) = \frac{6 + 4x + x^2}{6 - 2x}.$$

第4章

1. 1) $A_{-1} = A_1 = h/3$, $A_0 = 4h/3$, 具有 3 次代数精度, 2) $A_{-1} = A_1 = 8h/3$, $A_0 = -4h/3$, 3 次代数精度, 3) $x_1 = -0.28990$, $x_2 = 0.62660$ 或 $x_1 = 0.68990$, $x_2 = -0.12660$, 2 次代数精度, 4) $= 1/12$, 3 次代数精度; 2. 1) $T_8 = 0.11140$, $S_4 = 0.11157$, 2) $T_{10} = 1.39148$, $S_5 = 1.45471$, 3) $T_4 = 17.22774$, $S_2 = 17.32222$, 4) $T_6 = 1.03562$, $S_3 = 1.03577$; 4. $S_l = 0.63233$, 误差 0.00035;

6. n=213, 用辛普森公式区间分为 8 等分; 8. (1) 0.71327, (2) 0, (3) 10.1517434; 9. n=2, I=10.9484, n=3, I=10.95014, (准确值 I=10.9517032); 10. 48708 km; 11. 3.14158; 12. 1) 1.09863, 2) 1.08940, 1.09862, 3) 1.09854; 13. 用三点公式求得一阶导数值 -0.247, -0.217, -0.189, 用五点公式得 -0.2483, -0.2163, -0.1883.

第 5 章

4 .

$$\begin{aligned}
 l_{ii} &= a_{ii} \quad (i = 1, 2, \dots, n), \\
 u_{ij} &= a_{ij}/l_{ii} \quad (j = 2, 3, \dots, n), \\
 l_{ik} &= a_{ik} - \sum_{r=1}^{k-1} l_{ir} u_{rk} \quad (i = k, \dots, n), \\
 u_{kj} &= a_{kj} - \sum_{r=1}^{k-1} l_{kr} u_{rj} / l_{kk} \quad (j = k+1, \dots, n);
 \end{aligned}$$

5 . (a) 设 \mathbf{U} 为上三角阵

$$\begin{aligned}
 x_n &= b_n / u_{nn}, \\
 x_i &= b_i - \sum_{j=i+1}^n u_{ij} x_j / u_{ii} \quad (i = n-1, n-2, \dots, 1),
 \end{aligned}$$

(b) $n(n+1)/2$,(c) 记 \mathbf{U}^{-1} 的元素为 s_{ij} , \mathbf{U} 的元素记为 u_{ij} :

$$\begin{aligned}
 s_{ii} &= 1/u_{ii} \quad (i = 1, 2, \dots, n), \\
 s_{ij} &= -\sum_{k=i+1}^j u_{ik} s_{kj} / u_{ii},
 \end{aligned}$$

$i = n-1, n-2, \dots, 1, j = i+1, \dots, n;$

7 .

$$\mathbf{A}^{-1} = \begin{pmatrix} -0.04705885 & 0.5882353 & -0.2705882 & -0.9411764 \\ 0.3882353 & -0.3529412 & 0.4823529 & 0.7647058 \\ -0.2235294 & 0.2941177 & -0.03529412 & -0.4705882 \\ -0.03529412 & -0.05882353 & 0.04705882 & 0.2941176 \end{pmatrix}$$

8 . (1) $\mathbf{x}_1 = -1/2$, $\mathbf{x}_2 = -2/3$, $\mathbf{x}_3 = -3/4$, $\mathbf{x}_4 = -4/5$, (2) 解 $\mathbf{Ly} = \mathbf{f}$, $\mathbf{y} = (1/2, 1/3, 1/4, 1/5, 1/6)^T$, (3) 解 $\mathbf{Ux} = \mathbf{y}$, $\mathbf{x} = (5/6, 2/3, 1/2, 1/3, 1/6)^T$;

9 . $\mathbf{x} = (1.1111, 0.77778, 2.55556)^T$;

10 . (a) \mathbf{A} 不能分解为三角阵的乘积, 但换行后可以.(b) \mathbf{B} 可以但不唯一, \mathbf{C} 可以且唯一;

11 . $\mathbf{A}_1 = 1.1$, $\mathbf{A}_{-1} = 0.8$, $\mathbf{A}_{-2} = 0.825$, $\mathbf{A}_{-F} = 0.8426$;18 . $\text{cond}(\mathbf{A})_1 = 39601$, $\text{cond}(\mathbf{A})_2 = 39206$.

第 6 章

1 . (a) 两种方法均收敛,(b) 用雅可比迭代法迭代 18 次,

$\mathbf{x}^{(18)} = (-3.9999964, 2.9999739, 1.9999999)^T$, 用高斯-塞德尔迭代法迭代 8 次,

$$\mathbf{x}^{(8)} = (-4.000036, 2.999985, 2.000003)^T;$$

2 . (a) 雅可比迭代法不收敛,高斯-塞德尔迭代法收敛,

(b) 雅可比迭代法收敛,高斯-塞德尔迭代法不收敛;

5 . $= 1.03$ 时迭代 5 次达到精度要求 $\mathbf{x}^{(5)} = (0.5000043, 0.1000001, -0.4999999)^T$, $= 1$ 时迭代 6 次达到精度要求 $\mathbf{x}^{(6)} = (0.5000038, 0.1000002, -0.4999995)^T$, $= 1.1$ 时迭代 6 次达到精度要求, $\mathbf{x}^{(6)} = (0.5000035, 0.9999989, -0.5000003)^T$;

6 . $= 0.9$, 迭代 8 次时达到精度要求,

$$\mathbf{x}^{(8)} = (-4.000027, 0.2999989, 0.2000003)^T;$$

第 7 章

1 . 0 512;

2 . 1) 和 2) 收敛, 3) 发散;

3 . 1) 二分 14 次得 0.0905456, 2) 迭代 5 次得 0.0905264;

5 . (2) $x_0 = 1.5$, $x_2 = 1.465572$, (3) $x_0 = 1.5$, $x_3 = 1.465571$;

7 . (1) $x_3 = 1.8794$, (2) $x_3 = 1.8794; 8.4.49342$;

11 . 牛顿法 $x_{20} = 1.895494$ 其他方法迭代次数 $n=4, 13, 10, 723805$;

14 . $- (n-1)/2^n a, (n+1)/2^n a; 15.1/4 a$;

16 . $(1.526824, 0.508139)^T$.

第 8 章

1 . (a) 取 $\mathbf{v} = (1, 1, 1)^T$, $\mathbf{v}_1 = 9.6058$, $\mathbf{x} = (1, 0.6056, -0.3945)^T$,

(b) 取 $\mathbf{v} = (1, 1, 1)^T$, $\mathbf{v}_1 = 8.86951$, $\mathbf{x} = (-0.60422, 1, 0.15094)^T$; 2 . $= 7.288$, $\mathbf{x} = (1, 0.5229, 0.2422)^T$;

$$5. \mathbf{u} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & -\frac{3}{5} & -\frac{4}{5} \\ 0 & -\frac{4}{5} & \frac{3}{5} \end{pmatrix}, \mathbf{A}_2 = \begin{pmatrix} 1 & -5 & 0 \\ -5 & \frac{77}{25} & \frac{14}{25} \\ 0 & \frac{14}{25} & -\frac{23}{25} \end{pmatrix};$$

7. (a) \mathbf{A} 的特征值为 $\lambda_1 = \frac{1}{2} + \frac{33}{2}$, $\lambda_2 = 2$, $\lambda_3 = \frac{1}{2} - \frac{33}{2}$,

(b) \mathbf{B} 的特征值为 $\lambda_1 = 2 + 3$, $\lambda_2 = 2$, $\lambda_3 = 2 - 3$.

选取位移 $s_k = b_{33}^{(k)}$,

$$\mathbf{B}_s = \begin{pmatrix} 3.7316925974 & 0.0249060210 & 0.0 \\ 0.0249060210 & 2.0003582102 & 0.0 \\ 0.0 & 0.2679491924 & 0.0 \end{pmatrix}$$

其中 $| | < 5 \times 10^{-11}$;

$$8. \mathbf{Q} = \frac{1}{3} \begin{pmatrix} 1 & 2 & 2 \\ 2 & 1 & -2 \\ 2 & -2 & 1 \end{pmatrix}, \mathbf{R} = \begin{pmatrix} 3 & -3 & 3 \\ 3 & -3 & 3 \\ 3 & 3 & 3 \end{pmatrix};$$

9. \mathbf{A} 的特征值 2 ± 6 , \mathbf{A}^{-1} 特征值 $\frac{1}{6} \pm \frac{1}{2}$.

第 9 章

1. 0, 0.0010, 0.0050; 2. 0.145; 4. 0.500, 1.142, 2.501, 7.245;

5. 1) 1.2428, 1.5836, 2.0442, 2.6510, 3.4365,

2) 1.7276, 2.7430, 4.0942, 5.8292, 7.9960;

8. (1) $0 < h < 0.02$, (2) $0 < h < 0.0278$, (3) $0 < h < +\infty$;

9. 显式 0.626, 隐式 0.633, 真值 0.6321;

10. $-\frac{5}{8}h^3y(x_n)$; 12. 刚性比 $s=20$, $0 < h < 0.139$.