

■ 高等学校理工科数学类规划教材

数值分析

NUMERICAL ANALYSIS

(第二版)

主编 王金铭 副主编 谢彦红 杜洪波



大连理工大学出版社
DALIAN UNIVERSITY OF TECHNOLOGY PRESS

■ 高等学校理工科数学类规划教材

数 值 分 析

(第二版)

NUMERICAL ANALYSIS

主 编 王金铭

副主编 谢彦红 杜洪波



大连理工大学出版社
DALIAN UNIVERSITY OF TECHNOLOGY PRESS

图书在版编目(CIP)数据

数值分析 / 王金铭主编. —2 版. —大连:大连理工大学出版社,
2010.8

高等学校理工科数学类规划教材

ISBN 978-7-5611-3753-6

I . 数… II . 王… III . 数值计算—高等学校—教材
IV . O241

中国版本图书馆 CIP 数据核字(2007)第 129642 号

大连理工大学出版社出版

地址:大连市软件园路 80 号 邮政编码:116023

发行:0411-84708842 传真:0411-84701466 邮购:0411-84703636

E-mail:dutp@dutp.cn URL:<http://www.dutp.cn>

大连理工印刷有限公司印刷 大连理工大学出版社发行

幅面尺寸:185mm×260mm 印张:12.5 字数:289 千字

2007 年 8 月第 1 版

2010 年 8 月第 2 版

2010 年 8 月第 3 次印刷

责任编辑:于建辉 王伟 责任校对:捷琳

封面设计:宋蕾

ISBN 978-7-5611-3753-6

定 价:25.00 元

前　言

随着计算机技术的快速发展和普及,科学计算已经成为继实验与理论后的第三种科研方法,与科学计算密切相关的数值分析课程也已经成为高等学校理工科的一门重要课程。

本书着重介绍与科学计算有关的数值分析的基本方法,在强调基本概念和理论阐释的基础上非常重视实际应用,特别是数值方法在计算机上的实现。本书在理论分析方面力求完整的前提下,适当减少抽象的理论叙述,加强算法与实际计算能力的培养,特别注重分析数值方法的构造思想。此外,本书还适当介绍了一些数值方法上的最新成果,如(循环)块三对角线性方程组的求解方法、预处理共轭梯度法、多重网格法等;同时每章给出了两个典型方法的 C 语言程序供读者参考。

本书共分 10 章,分别介绍了数值分析中常用的数值方法和建立数值方法的基本原理。第 1 章绪论部分介绍了数值分析的研究对象与特点,误差的来源、分类及度量,误差稳定性分析与防止误差的原则。第 2~5 章为数值代数的基本内容。第 2 章介绍了线性方程组的直接法,主要包括高斯消去法、高斯列主元消去法、高斯-若当消去法、直接三角分解法及特殊线性方程组的直接三角分解法等;第 3 章介绍了线性方程组的迭代法,主要包括向量与矩阵范数、线性方程组的误差分析、三种常见的简单迭代法(雅克比迭代法、高斯-赛德尔迭代法、超松弛迭代法)、共轭梯度法及预处理共轭梯度法等;第 4 章介绍了非线性方程与方程组的数值解法,主要包括非线性方程的迭代法及其收敛性与收敛阶、牛顿法及其变形、非线性方程组的牛顿法及拟牛顿法等;第 5 章介绍了矩阵特征值问题计算,主要包括幂法与反幂法、Jacobi 方法、QR 方法等。第 6~8 章为数值逼近的基本内容。第 6 章介绍了函数的插值法,主要包括拉格朗日插值,差商型、差分型牛顿插值,埃尔米特插值,三次样条插值等;第 7 章介绍了最佳平方逼近及最小二乘法,主要包括连续函数的最佳平方逼近、离散函数的最小二乘法等;第 8 章介绍了数值积分与数值微分,主要包括插值型求积公式、等距节点的求积公式、龙贝格算法、高斯求积公式、重积分的计算公式、数值微分公式等。第 9、10 章为常微分方程数值解法的基本内容。第 9 章介绍了常微分方程初值问题的数值解法,主要包括欧拉法、改进欧拉法、四阶龙格-库塔法、线性多步法、一阶方程组与高阶方程的数值解法等;第 10 章介绍了常微分方程边值问题的数值解法,主要包括求解二阶常微分方程第一边值问题的打靶法、有限差分法及多重网格法等。

本次修订从以下四个方面进行：

- (1)对第2~4章的部分内容作了顺序调整；
- (2)对第2~4、7章增加了一些例题，同时对习题作了部分修改；
- (3)对每章的程序作了认真修改、逐一调试，通过算例进行验证；
- (4)针对书中存在的一些不妥之处和错误，对全书进行整体加工，统一协调。

本书第一版的第1~4、9、10章由王金铭编写；第5、6章由刘艳秋编写；第7、8章由陈欣编写。第二版的第1~5、9、10章由沈阳工业大学王金铭修订，第7、8章由沈阳化工学院谢彦红修订，第6章及各章程序由沈阳工业大学杜洪波修订。全书由王金铭统稿并最后定稿。

在本书的出版过程中，得到了沈阳工业大学理学院、研究生学院、教务处的大力支持与帮助，在此一并表示真诚的感谢；同时对审阅本书的专家表示感谢，谢谢你们所提出的宝贵意见。

本书可作为本科生和工科研究生数值分析课程的教材或教学参考书，也可供科技工作者和工程技术人员学习、参考。

由于编者水平有限，书中错误和不足之处在所难免，恳请读者批评指正。可通过以下的E-mail与作者联系：

wangjm_2004@163.com, jcjf@dutp.cn。

编者

2010年7月

目 录

第1章 绪 论 / 1

1.1 数值分析的概念与特点 / 1

 1.1.1 数值分析的概念 / 1

 1.1.2 数值分析的特点 / 1

1.2 误 差 / 2

 1.2.1 误差来源与分类 / 2

 1.2.2 误差的度量 / 2

1.3 数值稳定性与避免误差危害 / 3

 1.3.1 算法的数值稳定性 / 3

 1.3.2 避免误差危害的原则 / 5

习题 1 / 6

第2章 解线性方程组的直接法 / 7

2.1 高斯消去法 / 7

 2.1.1 上三角形方程组求解 / 7

 2.1.2 高斯消去法的基本思想 / 8

 2.1.3 解 n 阶线性方程组的高斯消去法 / 8

 2.1.4 矩阵的三角分解 / 10

 2.1.5 高斯消去法的计算量 / 11

2.2 高斯主元素消去法 / 12

 2.2.1 高斯列主元消去法 / 13

 2.2.2 高斯-若当消去法 / 13

2.3 高斯消去法的变形 / 15

 2.3.1 直接三角分解法 / 15

 2.3.2 特殊矩阵的直接三角分解 / 18

 2.3.3 列主元三角分解法 / 23

本章典型方法的 C 语言程序 / 25

习题 2 / 27

第3章 解线性方程组的迭代法 / 29

3.1 向量和矩阵的范数 / 29

 3.1.1 向量的数量积及其性质 / 29

 3.1.2 向量范数 / 30

 3.1.3 矩阵范数 / 31

 3.1.4 线性方程组的摄动分析 / 33

3.2 简单迭代法 / 34

 3.2.1 迭代法的基本思想 / 34

 3.2.2 简单迭代法的构造及相关概念 / 35

3.2.3 三种常见的简单迭代法 / 35

3.3 简单迭代法的收敛性 / 41

 3.3.1 迭代法收敛的基本定理 / 41

 3.3.2 迭代法收敛的误差估计 / 42

 3.3.3 三种常见的简单迭代法的简单判别方法 / 43

3.4 共轭梯度法 / 46

 3.4.1 与线性方程组等价的变分问题 / 46

 3.4.2 最速下降法 / 47

 3.4.3 共轭梯度法 / 49

 * 3.4.4 预处理共轭梯度法 / 52

本章典型方法的 C 语言程序 / 53

习题 3 / 55

第4章 非线性方程(组)的数值解法 / 58

4.1 引 言 / 58

4.2 二分法 / 59

4.3 迭代法 / 60

 4.3.1 迭代格式的构造 / 60

 4.3.2 迭代法的几何解释 / 60

 4.3.3 计算步骤 / 60

 4.3.4 收敛性与误差估计 / 61

 4.3.5 局部收敛性 / 63

 4.3.6 迭代法的收敛阶 / 63

 4.3.7 迭代收敛的加速方法 / 64

4.4 牛顿迭代法 / 66

 4.4.1 一般牛顿法 / 66

 4.4.2 牛顿法的变形 / 68

 * 4.5 解非线性方程组的牛顿迭代法 / 71

 4.5.1 Newton 法 / 71

 4.5.2 拟 Newton 法 / 72

本章典型方法的 C 语言程序 / 73

习题 4 / 75

第5章 矩阵特征值问题 / 77

5.1 幂法与反幂法 / 77

 5.1.1 幂 法 / 78

 5.1.2 反幂法 / 81

5.2 计算实对称矩阵特征值的雅可比方法 / 82

* 5.3 QR 方法简介 / 86

 5.3.1 矩阵 A 的 QR 分解 / 86

 5.3.2 QR 方法 / 87

本章典型方法的 C 语言程序 / 87

习题 5 / 90

第 6 章 插值法 / 91

6.1 问题的提出 / 91

 6.1.1 插值函数的概念 / 91

 6.1.2 插值多项式的存在唯一性 / 92

6.2 拉格朗日插值多项式 / 92

 6.2.1 线性插值和抛物插值 / 93

 6.2.2 拉格朗日插值多项式 / 95

 6.2.3 插值余项 / 96

6.3 差商、差分及牛顿插值公式 / 99

 6.3.1 差商及牛顿插值公式 / 99

 6.3.2 差分及等距节点牛顿插值公式 / 101

6.4 埃尔米特插值 / 104

6.5 分段低次插值 / 106

 6.5.1 高次插值的误差分析 / 106

 6.5.2 分段低次插值 / 107

6.6 三次样条插值 / 109

 6.6.1 三次样条插值函数 / 109

 6.6.2 三弯矩方法 / 110

本章典型方法的 C 语言程序 / 112

习题 6 / 114

第 7 章 最佳平方逼近及最小二乘法 / 116

7.1 函数的内积与正交多项式 / 116

 7.1.1 函数的内积及其性质 / 116

 7.1.2 正交多项式 / 117

 7.1.3 勒让德多项式 / 118

7.2 最佳平方逼近多项式 / 118

 7.2.1 基本概念及其计算 / 118

 7.2.2 用勒让德多项式作最佳平方逼近 / 120

7.3 最小二乘法 / 121

 7.3.1 最小二乘问题 / 121

 7.3.2 用最小二乘法求数据的拟合曲线 / 122

 7.3.3 用正交多项式作最小二乘拟合 / 125

7.3.4 利用最小二乘方法解超定方程组 / 126

本章典型方法的 C 语言程序 / 128

习题 7 / 130

第 8 章 数值积分与数值微分 / 131

8.1 数值积分问题的提出 / 131

 8.1.1 插值型求积公式 / 131

 8.1.2 插值型求积公式的截断误差与代数精度的概念 / 132

8.2 等距节点的求积公式 / 133

 8.2.1 柯特斯系数 / 133

 8.2.2 几种低阶牛顿-柯特斯公式
的截断误差 / 135

 8.2.3 复化求积公式与截断误差 / 136

8.3 变步长求积算法 / 138

 8.3.1 变步长梯形求积算法 / 138

 8.3.2 龙贝格算法 / 138

8.4 高斯求积公式 / 141

 8.4.1 一般理论 / 141

 8.4.2 高斯-勒让德求积公式 / 143

8.5 重积分的近似计算 / 145

8.6 数值微分 / 147

 8.6.1 数值微分问题的提出 / 147

 8.6.2 插值型求导公式及截断误差 / 148

本章典型方法的 C 语言程序 / 149

习题 8 / 151

第 9 章 常微分方程初值问题的数值解法 / 153

9.1 问题的提出 / 153

9.2 欧拉方法 / 154

 9.2.1 欧拉公式 / 154

 9.2.2 后退欧拉公式 / 155

 9.2.3 改进欧拉公式 / 156

 9.2.4 欧拉两步公式 / 158

9.3 龙格-库塔方法 / 160

 9.3.1 龙格-库塔方法的基本思想 / 160

 9.3.2 二阶龙格-库塔公式 / 160

 9.3.3 高阶龙格-库塔公式 / 161

 9.3.4 变步长的龙格-库塔方法 / 163

9.4 线性多步法 / 164

 9.4.1 基于数值积分的构造方法 / 164

 9.4.2 阿当姆斯内插公式 / 165

9.4.3 阿当姆斯外推公式及其阿当姆斯 预测-校正系统 / 166	10.3.2 多重网格法 / 179
9.5 一阶方程组与高阶方程 / 168	本章典型方法的 C 语言程序 / 180
9.5.1 一阶方程组 / 168	习题 10 / 182
9.5.2 化高阶方程为一阶方程组 / 169	参考答案与提示 / 183
本章典型方法的 C 语言程序 / 171	习题 1 / 183
习题 9 / 172	习题 2 / 183
第 10 章 常微分方程边值问题的数值解法 / 174	习题 3 / 183
10.1 打靶法 / 174	习题 4 / 184
10.2 有限差分法 / 175	习题 5 / 184
10.2.1 解二阶线性常微分方程第一边值 问题的差分方法 / 176	习题 6 / 185
10.2.2 解二阶非线性常微分方程第一边值 问题的差分方法 / 177	习题 7 / 186
10.3 多重网格法 / 177	习题 8 / 186
10.3.1 二重网格法 / 177	习题 9 / 187
	习题 10 / 189
	参考文献 / 190

第1章 絮 论

1.1 数值分析的概念与特点

1.1.1 数值分析的概念

数值分析是研究适合于计算机上使用的求解各种数学问题的数值计算方法及与此相关的理论的一门数学课程.

数值分析是一门内容丰富,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实验的高度技术性的特点.其内容包括线性方程组的数值解法、非线性方程(组)数值解法、矩阵特征值计算方法、函数的数值逼近、数值积分与数值微分、常微分方程数值解法等.

1.1.2 数值分析的特点

数值分析是一门数学课程,但它不像纯数学那样研究数学本身的理论,而是把数学理论与计算机紧密地结合起来,是一门与计算机密切相关的实用性很强的学科.

数值分析的特点概括起来可分为:

(1) 面向计算机

要根据计算机的特点,对数学问题提出或选择实际可行的有效算法.

(2) 算法应具有可靠的误差分析

由于计算机只能近似地表示实数,且任一算法只能在有限的时间内通过有限次运算完成,因此算法的收敛性和数值稳定性应得到保证,算法引起的误差应得到有效的控制.这些问题的解决往往需要建立相应的数学理论基础.

(3) 要有好的计算复杂性

计算复杂性问题是数值计算关心的一个重要问题,主要包括时间复杂性与空间复杂性.时间复杂性是指算法在有限的时间内结束运算,且所用时间尽可能少.空间复杂性是指算法所需的计算机的内存量不能太大,且所需存贮空间尽可能小.

(4) 要有数值实验

(5) 数值分析中有些方法虽然在数学理论上尚不严格和完善,但通过实际计算、对比分析等手段证明是行之有效的方法,也常常采用.

任何一个算法除了在理论上说明或论证其收敛性、数值稳定性外,还应通过计算机的具体数值试验来证明算法的可行性和有效性.

1.2 误 差

1.2.1 误差来源与分类

用计算机数值求解实际问题时经历以下过程:

实际问题 → 数学模型 → 数值计算方法 → 程序设计 → 上机运算 → 计算结果
按照上述过程可将误差分为以下四种.

(1) 模型误差

数学模型是对实际问题进行抽象、简化而得到的,其与实际问题之间的误差称为模型误差.由于模型误差不属于数值分析课程研究范畴,在这里不予讨论.

(2) 观测误差

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,这些物理量显然也包含误差.这种由观测产生的误差称为观测误差.在这里也不予讨论.

(3) 截断误差或方法误差

当数学模型不能得到精确解时,通常要用数值方法求其近似解,其近似解与精确解之间的误差称为截断误差.例如用

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

近似代替 $f(x)$,其余项

$$R_n(x) = f(x) - p_n(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!}x^{n+1}$$

称为截断误差.

(4) 舍入误差

由于计算机的字长有限,原始数据在计算机上表示及计算过程中可能产生误差,这种误差主要由舍入产生,所以称为舍入误差.例如,用 3.141 59 近似代替 π ,产生的误差

$$R = \pi - 3.141 59 = 0.000 002 6 \dots$$

就是舍入误差.

1.2.2 误差的度量

1. 绝对误差与绝对误差限

设 x^* 为准确值, x 为近似值,称

$$e = x^* - x$$

为绝对误差,简称误差.通常准确值 x^* 是未知的,因此误差 e 的准确值也是难以得到的.

但是根据测量工具的精密度或计算的实际情况可以确定误差的绝对值不超过某一正数 ϵ , ϵ 称为绝对误差限, 简称误差限. 精确值 x^* 、近似值 x 及误差限 ϵ 三者的关系为

$$x - \epsilon \leqslant x^* \leqslant x + \epsilon$$

通常记为

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

2. 相对误差与相对误差限

由于绝对误差没有反映出其相对于精确值的大小或占精确值的比例, 有时不能很好地度量近似值的精确程度, 因此通常使用相对误差. 称

$$\epsilon_r = \frac{\epsilon}{x^*} \approx \frac{\epsilon}{x}$$

为近似值 x 的相对误差. 类似绝对误差限的定义也可以确定相对误差的绝对值不超过某一正数 ϵ_r , 称它为近似值 x 的相对误差限.

3. 有效数字

设 x 是 x^* 的一个近似值, 写成

$$x = \pm 10^k \times [a_1 + a_2 \times 10^{-1} + \cdots + a_n \times 10^{-(n-1)}]$$

其中 $a_i (i = 1, 2, \dots, n)$ 是 $0, 1, \dots, 9$ 中的一个数字, $a_1 \neq 0$, k 为整数. 如果

$$|x - x^*| \leqslant 0.5 \times 10^{k+1-n}$$

则称 x 为 x^* 的具有 n 位有效数字的近似值.

例如, 用 3.14 近似 π , 有三位有效数字, 用 3.1416 则有五位有效数字. 这里, 近似值的获得可能由准确值四舍五入得到, 也可以由各种近似计算方法得到.

近似值 x 的有效数字与相对误差限有如下关系:

(1) 若 x 具有 n 位有效数字, 则其相对误差限为

$$\epsilon_r \leqslant \frac{1}{2a_1} \times 10^{1-n}$$

(2) 若 x 的相对误差限为

$$\epsilon_r \leqslant \frac{1}{2(a_1 + 1)} \times 10^{1-n}$$

则 x 至少具有 n 位有效数字.

显然, 近似值的有效数位数越多, 相对误差限越小; 反之亦然.

1.3 数值稳定性与避免误差危害

1.3.1 算法的数值稳定性

【例 1-1】 计算 $I_n = e^{-1} \int_0^1 x^n e^x dx (n = 0, 1, \dots)$, 并估计误差.

由分部积分可得计算 I_n 的递推公式

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

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

若计算出 I_0 , 代入式(1-1), 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

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

并取 $k = 7$, 用 4 位有效数字计算, 则得 $e^{-1} \approx 0.3679$, 截断误差 $R_7 = |e^{-1} - 0.3679| \leqslant \frac{1}{8!} \leqslant \frac{1}{4} \times 10^{-4}$. 当初值取为 $I_0 \approx 0.6321 = \tilde{I}_0$ 时, 用式(1-1) 递推的计算公式为

$$\begin{cases} \tilde{I}_0 = 0.6321 \\ \tilde{I}_n = 1 - n\tilde{I}_{n-1} \quad (n = 1, 2, \dots) \end{cases} \quad (\text{A})$$

计算结果见表 1-1 的 \tilde{I}_n 列. 用 \tilde{I}_0 近似 I_0 产生的误差 $E_0 = I_0 - \tilde{I}_0$ 就是初始误差, 它对后面计算结果是有影响的.

从表中看到 \tilde{I}_8 出现负值, 这与一切 $I_n > 0$ 相矛盾. 实际上, 由积分估值得

$$\frac{e^{-1}}{n+1} = e^{-1} \left(\min_{0 \leq x \leq 1} e^x \right) \int_0^1 x^n dx < I_n < e^{-1} \left(\max_{0 \leq x \leq 1} e^x \right) \int_0^1 x^n dx = \frac{1}{n+1} \quad (1-2)$$

因此, 当 n 较大时, 用 \tilde{I}_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 \tilde{I}_0 有误差 $E_0 = I_0 - \tilde{I}_0$, 由此引起以后各步计算的误差 $E_n = I_n - \tilde{I}_n$ 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 \tilde{I}_n 有误差 E_0 , 则 \tilde{I}_n 就有 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 \tilde{I}_8 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. 由式(1-2), 取 $n = 9$, 取

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

我们粗略取 $I_9 \approx \frac{1}{2} \left(\frac{1}{10} + \frac{e^{-1}}{10} \right) = 0.0684 = I_9^*$, 然后将式(1-1) 倒过来算, 即由 I_9^* 算出 $I_8^*, I_7^*, \dots, I_0^*$, 公式为

$$\begin{cases} I_9^* = 0.0684 \\ I_{n-1}^* = \frac{1}{n} (1 - I_n^*) \quad (n = 9, 8, \dots, 1) \end{cases} \quad (\text{B})$$

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 $|E_n^*| = \frac{1}{n!} |E_n^*|$, E_n^* 比 E_n 缩小了 $n!$ 倍, 因此, 尽管 E_9^* 较大, 但由于误差逐步缩小, 故可用 I_n^* 近似 I_n , 它表明计算公式(B) 是数值稳定的.

表 1-1

n	\bar{I}_n (用式(A)算)	I_n^* (用式(B)算)	n	\bar{I}_n (用式(A)算)	I_n^* (用式(B)算)
0	0.632 1	0.632 1	5	0.148 0	0.145 5
1	0.367 9	0.367 9	6	0.112 0	0.126 8
2	0.264 2	0.264 3	7	0.216 0	0.112 1
3	0.207 4	0.207 3	8	-0.728 0	0.103 5
4	0.170 4	0.170 8	9	7.552	0.068 4

定义 1-1 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,否则称此算法为数值不稳定的.

1.3.2 避免误差危害的原则

(1) 要使用数值稳定的计算公式

(2) 要尽量避免两相近数相减

在数值计算中,两相近数相减会使有效数字严重损失. 例如,用四位有效数字计算 $A = 10^7(1 - \cos 2^\circ)$ 的值(0.6092×10^4). 由于 $\cos 2^\circ = 0.9994$, 直接计算得 $A = 0.6000 \times 10^4$, 只有一位有效数字. 若利用倍角公式计算, $A = 10^7 \times 2 \times (\sin 1^\circ)^2 = 10^7 \times 2 \times 0.01745^2 = 0.6090 \times 10^4$, 具有三位有效数字. 如果无法改变算法以避免两相近数相减, 则采用双倍字长计算以增加有效数字.

(3) 要尽量避免大数“吃掉”小数

在数值计算中参加运算的数有时数量级相差很大,如果不注意采取相应措施,在它们的加、减法运算中,绝对值很小的数往往被绝对值较大的数“吃掉”,不能发挥其作用,造成计算结果失真. 例如,在八位十进制计算机上计算 $A = 10^{12} + 1000 - 10^{12}$, 直接计算有 $A = 0$, 显然是错误的,这时小数 1000 被大数 10^{12} 所“吃掉”. 如果改变顺序 $A = 10^{12} - 10^{12} + 1000$, 再计算可得正确结果 $A = 1000$.

(4) 要尽量避免绝对值很小的数作除数

在计算过程中,用绝对值很小的数作除数会使商的数量级增加,假设 x 和 y 的近似值分别是 \tilde{x} 和 \tilde{y} ,则 $z = x/y$ 的近似值是 $\tilde{z} = \tilde{x}/\tilde{y}$. 此时 z 的绝对误差

$$|e(z)| = |z - \tilde{z}| = \left| \frac{(x - \tilde{x})\tilde{y} + \tilde{x}(\tilde{y} - y)}{y\tilde{y}} \right| \approx \frac{|\tilde{y}| |x - \tilde{x}| + |\tilde{x}| |y - \tilde{y}|}{|\tilde{y}|^2}$$

可见,当 $|\tilde{y}|$ 很小时, \tilde{z} 的绝对值可能很大. 此外,当商过大时,或者其数值超出计算机表示的范围引起“溢出”现象,或者作为一个大数它将吃掉参与运算的一些小数.

(5) 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.

例如,对给定的 x ,计算多项式

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

的值. 如果采用逐项计算然后相加的算法:

$$p_n(x) = u_n + u_{n-1} + \cdots + u_0, \quad u_k = a_k x^k \quad (k = 0, 1, \dots, n)$$

一共需做

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

次乘法和 n 次加法. 如果把 $p_n(x)$ 改写为

$$p_n(x) = \{ \cdots [(a_n x + a_{n-1}) x + a_{n-2}] x + \cdots \} x + a_0$$

采用如下算法(秦九韶算法)

$$\begin{cases} b_n = a_n \\ b_k = b_{k+1}x + a_k & (k = n-1, \dots, 1, 0) \\ p_n(x) = b_0 \end{cases}$$

则只需 n 次乘法和 n 次加法运算.

再如, 利用 $\ln(1+x) = \sum_{n=1}^{\infty} (-1)^{n+1} \frac{x^n}{n}$ 计算 $\ln 2$, 若要精确到 10^{-5} , 要计算十万项求和. 这一方面使计算量很大; 另一方面舍入误差的积累也十分严重. 如果改用级数

$$\ln \frac{1+x}{1-x} = 2 \left[x + \frac{x^3}{3} + \frac{x^5}{5} + \cdots + \frac{x^{2n+1}}{2n+1} + \cdots \right]$$

取 $x = 1/3$, 只需计算前 9 项, 截断误差便小于 10^{-5} .

习题 1

1. 将下列各数四舍五入成五位有效数字

$$x_1 = 3.25894, x_2 = 3.25896, x_3 = 4.382000, x_4 = 0.000789247$$

2. 序列 $\{y_n\}$ 满足递推关系: $y_n = 100y_{n-1} - 2$ ($n = 1, 2, \dots$), 若 $y_0 = \sqrt{3} \approx 1.7321$, 计算到 y_5 时误差有多大? 这个计算过程稳定吗?

3. 设 $f(x) = 8x^5 + 4x^3 - 9x + 1$, 用秦九韶算法求 $f(3)$.

第2章 解线性方程组的直接法

求解线性方程组问题不仅在工程技术中涉及到,而且计算方法其他分支的研究(如样条插值,解非线性方程组问题,求解偏微分方程的差分法及有限元法等)也往往归结为此类问题,因此这是一个应用相当广泛的课题,也是数值代数中的一个重要研究课题.

考虑线性方程组

$$Ax = b \quad (2-1)$$

其中

$$A = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{pmatrix} \in \mathbb{R}^{n \times n}, \quad x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^n, \quad b = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix} \in \mathbb{R}^n$$

且设 A 为 n 阶非奇异矩阵.

直接法就是在没有舍入误差的情况下,通过有限步四则运算可以求得方程组精确解的方法.但由于实际计算中舍入误差是客观存在的,因而使用此类方法也只能得到近似解.直接法适合于求解中小型线性方程组.

2.1 高斯消去法

2.1.1 上三角形方程组求解

考虑系数矩阵为上三角矩阵的方程组

$$\left\{ \begin{array}{l} r_{11}x_1 + r_{12}x_2 + \cdots + r_{1n}x_n = b_1 \\ r_{22}x_2 + \cdots + r_{2n}x_n = b_2 \\ \vdots \\ r_{nn}x_n = b_n \end{array} \right. \quad (2-2)$$

其系数矩阵 $R = (r_{ij})$ [$r_{ii} \neq 0$, $r_{ij} = 0$ ($i > j$)] 是非奇异上三角矩阵.

求解方程组(2-2)的计算公式为

$$x_n = \frac{b_n}{r_{nn}}, \quad x_k = \frac{b_k - \sum_{j=k+1}^n r_{kj}x_j}{r_{kk}} \quad (k = n-1, n-2, \dots, 1) \quad (2-3)$$

式(2-3)称为“回代”公式,用式(2-3)求式(2-2)的过程称为“回代”过程.

2.1.2 高斯消去法的基本思想

由于上三角形方程组求解非常容易,因此在求解一般线性方程组(2-1)时,通过使用初等行变换把它化成等价的上三角形方程组(2-2)来求解,这即为高斯消去法的基本思想.将一般线性方程组化成等价上三角形方程组的过程称为“消去”过程.

【例 2-1】 用高斯消去法解方程组

$$\begin{cases} 6x_1 - 2x_2 + 2x_3 + 4x_4 = 16 \\ 12x_1 - 8x_2 + 6x_3 + 10x_4 = 26 \\ 3x_1 - 13x_2 + 9x_3 + 3x_4 = -19 \\ -6x_1 + 4x_2 + x_3 - 18x_4 = -34 \end{cases} \quad (2-4)$$

解

$$(A : b) = \left[\begin{array}{ccccc} 6 & -2 & 2 & 4 & 16 \\ 12 & -8 & 6 & 10 & 26 \\ 3 & -13 & 9 & 3 & -19 \\ -6 & 4 & 1 & -18 & -34 \end{array} \right] \xrightarrow{\begin{array}{l} r_2 - 2r_1 \\ r_3 - \frac{1}{2}r_1 \\ r_4 + r_1 \end{array}} \left[\begin{array}{ccccc} 6 & -2 & 2 & 4 & 16 \\ 0 & -4 & 2 & 2 & -6 \\ 0 & -12 & 8 & 1 & -27 \\ 0 & 2 & 3 & -14 & -18 \end{array} \right] \\ \xrightarrow{\begin{array}{l} r_3 - 3r_2 \\ r_4 + \frac{1}{2}r_2 \end{array}} \left[\begin{array}{ccccc} 6 & -2 & 2 & 4 & 16 \\ 0 & -4 & 2 & 2 & -6 \\ 0 & 0 & 2 & -5 & -9 \\ 0 & 0 & 4 & -13 & -21 \end{array} \right] \xrightarrow{r_4 - 2r_3} \left[\begin{array}{ccccc} 6 & -2 & 2 & 4 & 16 \\ 0 & -4 & 2 & 2 & -6 \\ 0 & 0 & 2 & -5 & -9 \\ 0 & 0 & 0 & -3 & -3 \end{array} \right]$$

其中, r_i 表示增广矩阵 $(A : b)$ 的第 i 行. 利用公式(2-3)求得方程组的解为

$$x^* = (3, 1, -2, 1)^T$$

2.1.3 解 n 阶线性方程组的高斯消去法

将式(2-1)记为 $A^{(1)}x = b^{(1)}$, 其增广矩阵为

$$(A^{(1)} : b^{(1)}) = \left[\begin{array}{ccccc} a_{11}^{(1)} & a_{12}^{(1)} & \cdots & a_{1n}^{(1)} & b_1^{(1)} \\ a_{21}^{(1)} & a_{22}^{(1)} & \cdots & a_{2n}^{(1)} & b_2^{(1)} \\ \vdots & \vdots & & \vdots & \vdots \\ a_{n1}^{(1)} & a_{n2}^{(1)} & \cdots & a_{nn}^{(1)} & b_n^{(1)} \end{array} \right]$$

1. 消去过程

(1) 第 1 次消去

设 $a_{11}^{(1)} \neq 0$, 首先计算乘数 $l_{ii} = a_{ii}^{(1)} / a_{11}^{(1)}$ ($i = 2, 3, \dots, n$), 用 $-l_{ii}$ 乘 $(A^{(1)} : b^{(1)})$ 的第 1 行加到第 i 行 ($i = 2, 3, \dots, n$) 上, 得到新的增广矩阵

$$(A^{(2)} : b^{(2)}) = \left[\begin{array}{ccccc} a_{11}^{(1)} & a_{12}^{(1)} & \cdots & a_{1n}^{(1)} & b_1^{(1)} \\ 0 & a_{22}^{(2)} & \cdots & a_{2n}^{(2)} & b_2^{(2)} \\ \vdots & \vdots & & \vdots & \vdots \\ 0 & a_{n2}^{(2)} & \cdots & a_{nn}^{(2)} & b_n^{(2)} \end{array} \right] \quad (2-5)$$

其中, $A^{(2)}, b^{(2)}$ 的元素计算公式为

$$a_{ij}^{(2)} = a_{ij}^{(1)} - l_{i1} a_{j1}^{(1)}, \quad b_i^{(2)} = b_i^{(1)} - l_{i1} b_1^{(1)} \quad (i, j = 2, 3, \dots, n)$$

(2) 第 k ($1 \leq k \leq n-1$) 次消去

设第 1 次消去直至第 $k-1$ 次消去过程计算已经完成, 即已得到增广矩阵

$$(A^{(k)} : b^{(k)}) = \begin{pmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \cdots & a_{1k}^{(1)} & \cdots & a_{1n}^{(1)} & b_1^{(1)} \\ a_{22}^{(2)} & \cdots & a_{2k}^{(2)} & \cdots & a_{2n}^{(2)} & b_2^{(2)} & \\ \vdots & \vdots & & \vdots & & \vdots & \\ a_{kk}^{(k)} & \cdots & a_{kn}^{(k)} & b_k^{(k)} & & & \\ \vdots & & \vdots & & & & \\ a_{nk}^{(k)} & \cdots & a_{nn}^{(k)} & b_n^{(k)} & & & \end{pmatrix} \quad (2-6)$$

设 $a_{kk}^{(k)} \neq 0$, 计算乘数 $l_{ik} = a_{ik}^{(k)} / a_{kk}^{(k)}$ ($i = k+1, \dots, n$), 用 $-l_{ik}$ 乘 $(A^{(k)} : b^{(k)})$ 的第 k 行加到第 i 行 ($i = k+1, k+2, \dots, n$) 上, 得到新的增广矩阵 $(A^{(k+1)} : b^{(k+1)})$, 其中 $A^{(k+1)}$, $b^{(k+1)}$ 的元素计算公式为

$$a_{ij}^{(k+1)} = a_{ij}^{(k)} - l_{ik} a_{kj}^{(k)} \quad b_i^{(k+1)} = b_i^{(k)} - l_{ik} b_k^{(k)} \quad (i, j = k+1, \dots, n) \quad (2-7)$$

继续上述过程, 并设 $a_{kk}^{(k)} \neq 0$ ($k = 1, 2, \dots, n-1$) ($a_{kk}^{(k)}$ 称为主元素), 直到完成第 $n-1$ 次消去, 最后得到与方程组 (2-1) 等价的上三角形方程组的增广矩阵

$$(A^{(n)} : b^{(n)}) = \begin{pmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \cdots & a_{1n}^{(1)} & b_1^{(1)} \\ a_{22}^{(2)} & \cdots & a_{2n}^{(2)} & b_2^{(2)} & \\ \vdots & & \vdots & & \\ a_{nn}^{(n)} & b_n^{(n)} & & & \end{pmatrix}$$

其对应的上三角形方程组为

$$A^{(n)} x = b^{(n)} \quad (2-8)$$

2. 回代过程

利用式 (2-3) 求出式 (2-8) 的解.

在进行消去时, 如果遇到主元素为零, 例如 $a_{11}^{(1)} = 0$, 由于 A 为非奇异矩阵, 所以 A 的第一列一定有元素不等于零, 例如 $a_{i_1 1}^{(1)} \neq 0$, 于是可交换两行元素, 将 $a_{i_1 1}^{(1)}$ 调到 $(1, 1)$ 位置, 然后进行消去计算; 这时 $A^{(2)}$ 右下角矩阵 ($n-1$ 阶方阵) 亦为非奇异矩阵. 继续这个过程, 高斯消去法照样可以进行下去.

综上所述有以下定理.

定理 2-1 如果 A 为 n 阶非奇异矩阵, 则可通过高斯消去法 (及交换两行的初等变换) 将方程组 (2-1) 化为上三角形方程组 (2-8).

定理 2-2 消去过程中的主元素 $a_{kk}^{(k)}$ ($k = 1, 2, \dots, n-1$) 非零的充要条件是矩阵 A 的顺序主子式 $\Delta_k \neq 0$ ($k = 1, 2, \dots, n-1$), 即

$$\Delta_1 = a_{11} \neq 0, \quad \Delta_k = \begin{vmatrix} a_{11} & \cdots & a_{1k} \\ \vdots & & \vdots \\ a_{k1} & \cdots & a_{kk} \end{vmatrix} \neq 0 \quad (k = 2, 3, \dots, n-1)$$