

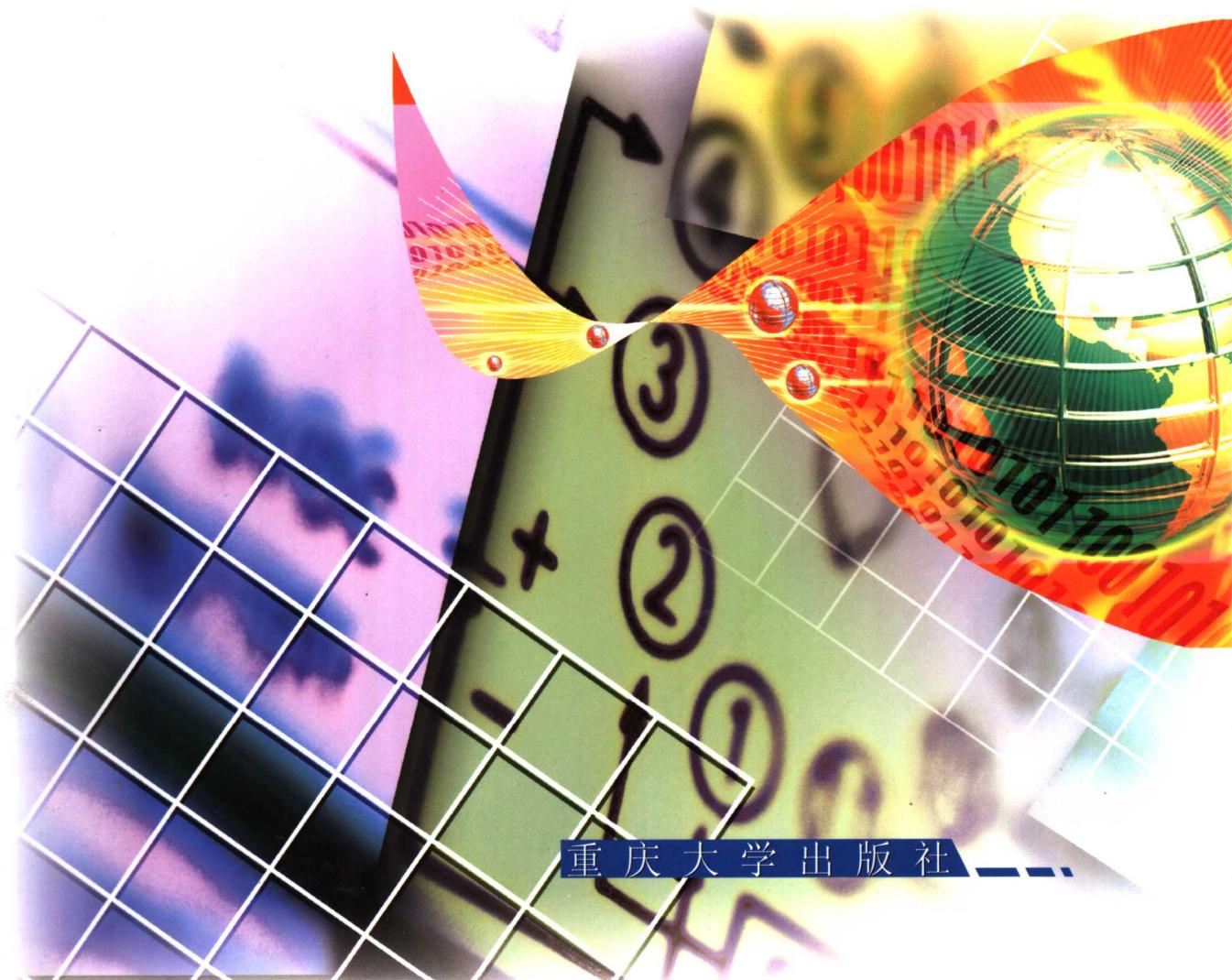
S

# 数值计算方法

SHUZHI JISUAN FANGFA

1 2 3 4 5 6 7 8 9 0  
1 2 3 4 5 6 7 8 9 0  
1 2 3 4 5 6 7 8 9 0  
1 2 3 4 5 6 7 8 9 0

- ◎ 主 编 郑继明  
◎ 副主编 刘 平 张清华



重庆大学出版社

# 数值计算方法

主 编 郑继明

副主编 刘 平 张清华

参 编 杨春德 刘 勇

重庆大学出版社

## 内 容 提 要

本书介绍了科学计算中最基本的数值计算方法。主要内容有：线性代数方程组的数值解法，非线性方程和方程组的迭代解法，矩阵特征值和特征向量的计算，函数的插值与曲线拟合，数值积分和常微分方程初值问题的数值解法。

本书可作高校理工科有关专业的教材，也可供有关科技人员参考。

### 图书在版编目(CIP)数据

数值计算方法/郑继明主编. —重庆:重庆大学出版社, 2005. 8

ISBN 7-5624-3438-7

I. 数... II. 郑... III. 数值计算—计算方法  
IV. 0241

中国版本图书馆 CIP 数据核字(2005)第 073692 号

### 数值计算方法

主 编 郑继明

副主编 刘 平 张清华

责任编辑:曾令维 邵孟春 版式设计:曾令维

责任校对:许 玲 责任印制:秦 梅

\*

重庆大学出版社出版发行

出版人:张鸽盛

社址:重庆市沙坪坝正街 174 号重庆大学(A 区)内

邮编:400030

电话:(023) 65102378 65105781

传真:(023) 65103686 65105565

网址:<http://www.cqup.com.cn>

邮箱:[fzk@cqup.com.cn](mailto:fzk@cqup.com.cn) (市场营销部)

全国新华书店经销

四川外语学院印刷厂印刷

\*

开本:787×1092 1/16 印张:8.75 字数:218 千

2005 年 8 月第 1 版 2005 年 8 月第 1 次印刷

印数:1—3 500

ISBN 7-5624-3438-7

定价:15.00 元

---

本书如有印刷、装订等质量问题,本社负责调换

版权所有,请勿擅自翻印和用本书

制作各类出版物及配套用书,违者必究。

# 前 言

在科学的研究和工程设计中经常需要做大量的数值计算。现在,数值分析方法与计算机技术相结合已深入到计算物理、计算力学、计算化学、计算生物学、计算经济学等各个领域,计算机上使用的数值计算方法已浩如烟海。本书是以理工科本科生为主要对象编写的,目的是使读者获得数值分析方法的基本概念,掌握适用于电子计算机的常用算法,具有基本的理论分析和实际计算能力。

本书只限于介绍科学计算中最基本的数值分析方法。主要内容有:线性代数方程组的数值解法,非线性方程和方程组的迭代解法,矩阵特征值和特征向量的计算,函数的插值与曲线拟合,数值积分和常微分方程初值问题的数值解法。

在学习数值分析时,要注意掌握数值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。此外,还要通过应用数值方法编程计算例子来提高使用各种数值方法解决实际问题的能力。

由于编者水平有限,书中难免有缺陷和疏漏,敬请广大读者批评指正。

编 者

2005 年 6 月

# 目 录

<b>第1章 数值计算中的误差</b>	1
1.1 引言	1
1.2 误差的种类及其来源	2
1.3 数值计算的误差	4
1.4 算法的数值稳定性	8
习 题1	11
<b>第2章 插值法</b>	13
2.1 插值问题	13
2.2 拉格朗日(lagrange)多项式插值	15
2.3 牛顿(Newton)插值	18
2.4 分段低次插值	21
2.5 样条插值	22
2.6 数值微分	28
习 题2	30
<b>第3章 曲线拟合的最小二乘法</b>	32
3.1 最小二乘法的提法	32
3.2 最小二乘法的求法	33
3.3 用正交多项式作最小二乘法	36
习 题3	37
<b>第4章 矩阵的特征值与特征向量</b>	39
4.1 乘幂法	39
4.2 乘幂法的加速方法	42
4.3 反幂法	44
4.4 雅可比(Jacobi)方法	45
4.5 QR方法	48
习 题4	55
<b>第5章 数值积分</b>	57
5.1 构造数值求积公式的基本方法	57
5.2 牛顿-科特斯求积公式	58
5.3 复化求积公式	63
5.4 龙贝格(Rumberg)求积算法	67
习 题5	69

<b>第6章 非线性方程及非线性方程组的解法</b>	70
6.1 二分法	70
6.2 迭代法	72
6.3 牛顿法	77
6.4 弦割法	79
6.5 解非线性方程组的迭代法	80
习题6	84
<b>第7章 解线性方程组的数值方法</b>	85
7.1 引言	85
7.2 高斯消去法	85
7.3 选主元素的高斯消去法	89
7.4 矩阵的三角分解	92
7.5 向量和矩阵的范数	100
7.6 解线性方程组的迭代法	103
7.7 病态方程组和迭代改善法	111
习题7	114
<b>第8章 常微分方程初值问题数值解法</b>	116
8.1 欧拉(Euler)方法	117
8.2 龙格-库塔(Runge-Kutta)方法	120
8.3 阿达姆斯(Adams)方法	122
8.4 收敛性与稳定性	124
8.5 方程组与高阶方程的数值解法	125
习题8	126
<b>附录 部分上机实习题</b>	128
<b>参考文献</b>	131

# 第 1 章

## 数值计算中的误差

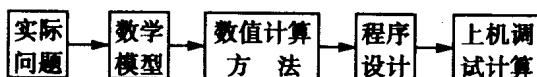
### 1.1 引言

随着计算机的发展和普及,数值计算已成为工程设计与科学的重要手段。掌握数值计算方法,会用计算机解决科学与工程实际中提出的数值计算问题,已成为科技人员必须具有的能力。

所谓数值计算方法,就是研究怎样利用计算尺、计算器、电子计算机等计算工具来求出数学问题的数值解,并对算法的收敛性、稳定性和误差进行分析、计算的全过程。它的理论与方法随计算工具的发展而发展。本书介绍计算机上常用的数值计算方法。

众所周知,传统的科学研究方法有两种:理论分析和科学实验。今天伴随着计算机技术的飞速发展和计算数学方法与理论的日益成熟,科学计算已成为第三种科学的研究的方法和手段。科学计算的物质基础是计算机,其理论基础是计算数学。随着科学技术的突飞猛进,无论是工农业生产还是国防尖端技术,例如机电产品的设计、建筑工程项目的设计、气象预报和新型尖端武器的研制、火箭的发射等,都有大量复杂的数值计算问题亟待解决。它们的复杂程度已达到非人工手算(包括使用计算器等简单的计算工具)所能解决的地步。数字电子计算机的出现和飞速发展大大推动了数值计算方法的发展,许多复杂的数值计算问题现在都可以通过使用计算机进行数值计算而得到妥善解决。

用数值计算的方法来解决工程实际和科学技术中的具体技术问题时,首先必须将具体问题抽象为数学问题,即建立起能描述并等价代替该实际问题的数学模型,例如各种微分方程、积分方程、代数方程等,然后选择合适的计算方法(算法),编制出计算机程序,最后上机调试并进行运算,以得出所欲求解的结果来。基本过程如下:



具体地说,数值计算方法首先要构造可计算出各种问题解的计算方法;然后分析方法的可靠性,即按此方法计算得到的解是否可靠,与精确解之差是否很小,以确保计算解的有效性;第

三,要分析方法的效率,分析比较求解同一问题的各种方法的计算速度和存储量,以便使用者根据各自的情况采用高效率的方法,节省人力、物力和时间,这样的分析是数值分析的一个重要部分。应当指出,数值方法的构造和分析是密切相关不可分割的。对于给定的数学问题,常常可以提出各种各样的数值计算方法。这里所说的“算法”,不仅是单纯的数学公式,而且是指由基本运算和运算顺序的规定所组成的整个解题方案和步骤。一般可以通过框图(流程图)来较直观地描述算法的全貌。

选定合适的算法是整个数值计算中非常重要的一环。例如,当计算多项式

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

的值时,若直接计算  $a_i x^i (i = 0, 1, \dots, n)$ ,再逐项相加,共需做

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

次乘法和  $n$  次加法。 $n = 10$  时需做 55 次乘法和 10 次加法。用著名的秦九韶(我国宋朝数学家)算法,将多项式  $P(x)$  改写成

$$P(x) = (((\dots((a_n x + a_{n-1})x + a_{n-2})x + \cdots + a_2)x + a_1)x + a_0)$$

来计算时,只要做  $n$  次乘法和  $n$  次加法即可。如当  $n = 10$  时,只要做 10 次乘法和 10 次加法。可见算法的优劣直接影响计算的速度和效率。

算法选得不恰当,不仅影响到计算的速度和效率,还会由于计算机计算的近似性和误差的传播、积累直接影响到计算结果的精度,有时甚至直接影响到计算的成败。不合适的算法会导致计算误差达到不能容许的地步,而使计算最终失败,这就是算法的数值稳定性问题。因此,最有效的算法,应当实用范围广,运算工作量少,需要存储单元少,逻辑结构简单,便于编写计算机程序,而且计算结果可靠。

## 1.2 误差的种类及其来源

数值计算过程中会出现各种误差,它们可分为两大类:一类是由于算题者在工作中的粗心大意而产生的,例如笔误将 886 误写成 868,以及误用公式等,这类误差称为“过失误差”或“疏忽误差”。它完全是人为造成的,只要在工作中仔细、谨慎,是完全可以避免的,我们就不再讨论它;而另一类为“非过失误差”,在数值计算中则往往是无法避免的,例如近似值带来的误差。在科学计算中误差来源一般有以下 4 个方面:模型误差、观测误差、截断误差和舍入误差等。

### 1.2.1 模型误差

在建模过程中,欲将复杂的物理现象抽象、归结为数学模型,往往只需忽略一些次要因素的影响,而对问题做某些必要的简化。这样建立起来的数学模型实际上必定只是所研究的复杂客观现象的一种近似的描述,它与真正客观存在的实际问题之间有一定的差别,这种误差称为模型误差。

### 1.2.2 观测误差

在建模和具体运算过程中所用的一些初始数据往往都是通过人们实际观察、测量得来的,

由于受到所用测量工具的限制或在数据的获取时受到随机因素的影响,这些数据都只能是近似的,即存在着误差,这种误差称为观测误差。

### 1.2.3 截断误差

在不少数值运算中常遇到超越计算,如微分、积分和无穷级数求和等,它们需用极限或无穷过程来求得。然而计算机却只能完成有限次算术运算和逻辑运算,因此需将解题过程化为一系列有限的算术运算和逻辑运算。这样就要对某种无穷过程进行“截断”。这种用有限过程代替无限过程所引起的误差,称为截断误差或方法误差。例如,函数  $\sin x$  和  $\ln(1+x)$  可分别展开为  $x$  的幂级数:

$$\begin{aligned}\sin x &= x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots \\ \ln(1+x) &= x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \dots\end{aligned}$$

若取级数的起始若干项的部分和作为  $x < 1$  时函数值的近似计算公式,例如取

$$\begin{aligned}\sin x &\approx x - \frac{x^3}{3!} + \frac{x^5}{5!} \\ \ln(1+x) &\approx x - \frac{x^2}{2} + \frac{x^3}{3}\end{aligned}$$

则由于它们的第四项和以后各项都舍弃了,自然产生了所谓的截断误差。

一般地,函数  $f(x)$  用泰勒(Taylor)多项式

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

近似代替,则截断误差是

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

例如:用多项式  $x - \frac{x^3}{3!} + \frac{x^5}{5!}$  计算  $\sin x$  的截断误差是

$$\left| R_4(x) \right| \leq \frac{x^7}{7!}$$

### 1.2.4 舍入误差

在数值计算过程中还会用到一些无穷小数,例如无理数和有理数中某些分数化出的无限循环小数,如

$$\pi = 3.14159265\dots$$

$$\sqrt{2} = 1.41421356\dots$$

$$\frac{1}{3!} = \frac{1}{6} = 0.166666\dots$$

等。而计算机受机器字长的限制,它所能表示的数据只能有一定的有限位数,这时就需把数据按四舍五入成一定位数的近似的有理数来代替。由此引起的误差称为舍入误差。

综上所述,数值计算中除了可以完全避免的过失误差外,还存在难以回避的模型误差、观测误差、截断误差和舍入误差。本课程主要考虑后两种误差。

## 1.3 数值计算的误差

### 1.3.1 绝对误差和相对误差

设某一个量的准确值(称为真值)为  $x$ , 其近似值为  $x^*$ , 则  $x$  与  $x^*$  的差

$$\varepsilon(x) = x - x^* \quad (1.1)$$

称为近似值  $x^*$  的绝对误差, 简称误差。

由于真值  $x$  往往是未知或无法知道的, 因此  $\varepsilon(x)$  的准确值也就无法求出。但一般可估计出此绝对误差  $\varepsilon(x)$  的上限, 也即可以求出一个正数  $\delta_1$ , 使

$$|\varepsilon(x)| = |x - x^*| \leq \delta_1 \quad (1.2)$$

$\delta_1$  称为近似值  $x^*$  的绝对误差限, 或称为精度。通常用

$$x = x^* \pm \delta_1$$

来表示近似值的精度。正数  $\delta_1$  越小, 表示该近似值  $x^*$  的精度越高。

在实际问题中, 判断一个近似值的精确度大小不仅要观察绝对误差大小, 还要考虑该量本身的大小。这就需要引进相对误差的概念。

相对误差定义为绝对误差与真值之比, 即

$$\varepsilon_r(x) = \frac{\varepsilon(x)}{x} = \frac{x - x^*}{x} \quad (x \neq 0) \quad (1.3)$$

例如测量 10 m 的长度时产生 1 cm 的误差与测量 1 m 的长度时产生 1 cm 的误差是大有区别的。虽然两者的绝对误差相同, 都是 1 cm, 但是前一种测量的相对误差为  $\frac{1}{1000}$ , 而后一种测量的相对误差则为  $\frac{1}{100}$ , 是前一种的十倍。

由式(1.3)可得

$$\varepsilon(x) = x \cdot \varepsilon_r(x) \quad (1.4)$$

相对误差不仅能表示出绝对误差来, 而且在估计近似值运算结果的误差时, 它比绝对误差更能反映出误差的特性。因此在误差分析中, 相对误差比绝对误差更为重要。

与绝对误差一样, 相对误差也无法准确求出, 但可以估计它的范围, 即可找到一个适当小的正数  $\delta_2$ , 称为近似值  $x^*$  的相对误差限, 即

$$|\varepsilon_r(x)| \leq \delta_2 \quad (1.5)$$

注: ① 相对误差没有量纲, 而绝对误差有量纲。

② 在实际计算中, 由于真值  $x$  总是无法知道的, 因此往往取

$$\varepsilon_r^*(x) = \frac{\varepsilon(x)}{x^*} \quad (1.6)$$

作为相对误差的另一定义。

③相对误差也常用百分数来表示：

$$\varepsilon_r(x) = \frac{\varepsilon(x)}{x} \times 100\%$$

这时称它为百分误差。

### 1.3.2 有效数字

在表示一个近似值时,为了同时反映其准确程度,常常用到“有效数字”的概念。例如对无穷小数或循环小数,可用四舍五入的办法来取其近似值。

**例 1.1** 我们知道, $\pi = 3.14159265\cdots$ 是一个无理数,按四舍五入考虑 $\pi$ 的不同近似值:

取一位数: $x_1^* = 3$ ,有 $| \pi - x_1^* | \leq 0.5 = \frac{1}{2}$ ;

取四位小数: $x_2^* = 3.1416$ ,有 $| \pi - x_2^* | \leq 0.00005 = \frac{1}{2} \times 10^{-4}$ ;

取五位小数: $x_3^* = 3.14159$ ,有 $| \pi - x_3^* | \leq 0.000005 = \frac{1}{2} \times 10^{-5}$ 。

这种近似值取法的特点是误差限为其末位数的半个单位。当近似值 $x^*$ 的绝对误差限是其某一位上的半个单位时,就称其“准确”到这一位,且从该位起直到前面第一位非零数字为止的所有数字都称为有效数字。

**定义** 设 $x$ 的近似值 $x^*$ 的规格化形式为

$$x^* = \pm 0.\alpha_1\alpha_2\cdots\alpha_n \times 10^m \quad (1.7)$$

其中, $\alpha_1, \alpha_2, \dots, \alpha_n$ 都是 $0 \sim 9$ 中的任一整数,且 $\alpha_1 \neq 0$ ; $n$ 是正整数, $m$ 是整数。若 $x^*$ 的误差限为

$$|\varepsilon(x)| = |x - x^*| \leq \frac{1}{2} \times 10^{m-n} \quad (1.8)$$

则称 $x^*$ 为具有 $n$ 位有效数字的有效数,或称它精确到 $10^{m-n}$ 。其中每一位数字 $\alpha_1, \alpha_2, \dots, \alpha_n$ 都是 $x^*$ 的有效数字。

**注:**①若式(1.7)中的 $x^*$ 是 $x$ 经四舍五入得到的近似值,则 $x^*$ 具有 $n$ 位有效数字。例如,3.1416是 $\pi$ 的具有五位有效数字的近似值,它精确到0.0001。

②有效数尾部的零不可随意省去,以免损失精度。

③另一种情况,例如 $x = 0.1524, x^* = 0.154$ 。这时 $x^*$ 的误差 $\varepsilon(x) = -0.0016$ ,其绝对值超过了0.0005(第三位小数的半个单位),但却没有超过0.005(第二位小数的半个单位),即

$$0.0005 < |x - x^*| \leq 0.005.$$

显然 $x^*$ 虽有三位小数但却只精确到第二位小数,因此它只具有二位有效数字。其中 $\alpha_1 = 1$ , $\alpha_2 = 5$ 都是准确数字,而第三位数字 $\alpha_3 = 4$ 就不再是准确数字了,就称它为存疑数字。

另外,由式(1.8)可知,从有效数字可以算出近似数的绝对误差限;有效数字的位数越多,其绝对误差限也就越小。不但如此,还可以从有效数字中求出其相对误差限。

当用式(1.7)表示的近似值 $x^*$ 具有 $n$ 位有效数字时,显然有

$$|x^*| \geq \alpha_1 \times 10^{m-1},$$

故由式(1.8)可知,其相对误差

$$\left| \varepsilon_r^*(x) \right| = \left| \frac{\varepsilon(x)}{x^*} \right| \leq \frac{\frac{1}{2} \times 10^{m-n}}{\alpha_1 \times 10^{m-1}} = \frac{1}{2\alpha_1} \times 10^{-n+1}.$$

故相对误差限为

$$\delta_2 = \frac{1}{2\alpha_1} \times 10^{-n+1} \quad (1.9)$$

式(1.9)说明  $x^*$  的有效数字位数越多,其相对误差限越小。由此可见,有效数字的位数反映了近似值的相对精确度。

**例 1.2** 当用 3.1416 来表示  $\pi$  的近似值时,它的相对误差限是多少?

**解** 3.1416 具有 5 位有效数字,  $\alpha_1 = 3$ , 由式(1.9)有

$$\delta_2 = \frac{1}{2\alpha_1} \times 10^{-n+1} = \frac{1}{2 \times 3} \times 10^{-5+1} = \frac{1}{6} \times 10^{-4}.$$

**例 1.3** 为了使  $x = \sqrt{20}$  的近似值  $x^*$  的相对误差小于 0.1%,问至少取几位有效数字?

**解** 因为  $\sqrt{20} = 4.47213\dots$ , 则近似值  $x^*$  中  $\alpha_1 = 4$ 。由式(1.9)知

$$\frac{1}{2 \times 4} \times 10^{-n+1} \leq 0.1\%$$

可解出  $n = 4$ 。即只要取 4 位有效数字,此时  $x^* = 4.472$  就能满足要求。

### 1.3.3 数值运算的误差

在实际的数值计算中,参与运算的数据往往都是些近似值,带有误差。这些数据误差在多次运算过程中会进行传播,使计算结果产生误差。而确定计算结果所能达到的精度显然是十分重要的,但这往往也是件很困难的事。不过,对计算误差做出一定的定量估计还是可以做到的。这里介绍一种常用的误差估计的一般公式,它是利用函数的泰勒(Taylor)展开得到的。

先从较简单的二元函数  $y = f(x_1, x_2)$  开始。设  $x_1^*$  和  $x_2^*$  分别是  $x_1$  和  $x_2$  的近似值,  $y^*$  是函数值  $y$  的近似值,且  $y^* = f(x_1^*, x_2^*)$ 。

函数  $f(x_1, x_2)$  在点  $(x_1^*, x_2^*)$  处的泰勒展开式为

$$f(x_1, x_2) = f(x_1^*, x_2^*) + \left[ \left( \frac{\partial f}{\partial x_1} \right)^* (x_1 - x_1^*) + \left( \frac{\partial f}{\partial x_2} \right)^* (x_2 - x_2^*) \right] + \\ \frac{1}{2!} \left[ \left( \frac{\partial^2 f}{\partial x_1^2} \right)^* (x_1 - x_1^*)^2 + 2 \left( \frac{\partial^2 f}{\partial x_1 \partial x_2} \right)^* (x_1 - x_1^*)(x_2 - x_2^*) + \left( \frac{\partial^2 f}{\partial x_2^2} \right)^* (x_2 - x_2^*)^2 \right] + \dots$$

式中,  $(x_1 - x_1^*) = \varepsilon(x_1)$  和  $(x_2 - x_2^*) = \varepsilon(x_2)$  一般都是小量值,如忽略高阶小量,即高阶的  $(x_1 - x_1^*)^2$  和  $(x_2 - x_2^*)^2$ , 则上式可简化为

$$f(x_1, x_2) \approx f(x_1^*, x_2^*) + \left( \frac{\partial f}{\partial x_1} \right)^* \varepsilon(x_1) + \left( \frac{\partial f}{\partial x_2} \right)^* \varepsilon(x_2)$$

因此  $y^*$  的绝对误差差

$$\varepsilon(y) = y - y^* = f(x_1, x_2) - f(x_1^*, x_2^*) \approx \left( \frac{\partial f}{\partial x_1} \right)^* \varepsilon(x_1) + \left( \frac{\partial f}{\partial x_2} \right)^* \varepsilon(x_2) \quad (1.10)$$

式中,  $\varepsilon(x_1)$  和  $\varepsilon(x_2)$  前面的系数  $\left(\frac{\partial f}{\partial x_1}\right)^*$  和  $\left(\frac{\partial f}{\partial x_2}\right)^*$  分别是  $x_1^*$  和  $x_2^*$  对  $y^*$  的绝对误差增长因子, 它们分别表示绝对误差  $\varepsilon(x_1)$  和  $\varepsilon(x_2)$  经过传播后增大或缩小的倍数。

由式(1.10)可得出  $y^*$  的相对误差

$$\varepsilon_r(y) = \frac{\varepsilon(y)}{y^*} \approx \left(\frac{\partial f}{\partial x_1}\right)^* \frac{\varepsilon(x_1)}{y^*} + \left(\frac{\partial f}{\partial x_2}\right)^* \frac{\varepsilon(x_2)}{y^*} = \frac{x_1^*}{y^*} \left(\frac{\partial f}{\partial x_1}\right)^* \varepsilon_r(x_1) + \frac{x_2^*}{y^*} \left(\frac{\partial f}{\partial x_2}\right)^* \varepsilon_r(x_2) \quad (1.11)$$

式中,  $\varepsilon_r(x_1)$  和  $\varepsilon_r(x_2)$  前面的系数  $\frac{x_1^*}{y^*} \left(\frac{\partial f}{\partial x_1}\right)^*$  和  $\frac{x_2^*}{y^*} \left(\frac{\partial f}{\partial x_2}\right)^*$  分别是  $x_1^*$  和  $x_2^*$  对  $y^*$  的相对误差增长因子, 它们分别表示相对误差  $\varepsilon_r(x_1)$  和  $\varepsilon_r(x_2)$  经过传播后增大或缩小的倍数。

由式(1.10)和式(1.11)两式可得加、减法运算的误差公式

$$\begin{aligned} \varepsilon(x_1 + x_2) &\approx \varepsilon(x_1) + \varepsilon(x_2), \varepsilon_r(x_1 + x_2) \approx \frac{x_1^*}{x_1 + x_2} \varepsilon_r(x_1) + \frac{x_2^*}{x_1 + x_2} \varepsilon_r(x_2) \\ \varepsilon(x_1 - x_2) &\approx \varepsilon(x_1) - \varepsilon(x_2), \varepsilon_r(x_1 - x_2) \approx \frac{x_1^*}{x_1^* - x_2} \varepsilon_r(x_1) - \frac{x_2^*}{x_1^* - x_2} \varepsilon_r(x_2) \end{aligned} \quad (1.12)$$

因此, 当  $x_1 \approx x_2$ , 即大小接近的两个同号近似值相减时, 由式(1.12)的第二式可知,  $|\varepsilon_r(x_1 - x_2)|$  可能会很大, 说明计算结果的有效数字将严重丢失, 计算精度很低。因此在实际计算中, 应尽量设法避开相近数的相减。当实在无法避免时, 可用变换计算公式的方法来解决。例如, 当要求计算  $\sqrt{3.01} - \sqrt{3}$ , 结果精确到第5位数字时, 至少取到  $\sqrt{3.01} = 1.7349352$  和  $\sqrt{3} = 1.7320508$ , 这样  $\sqrt{3.01} - \sqrt{3} = 2.8844 \times 10^{-3}$  才能达到具有5位有效数字的要求。如果变换算式:

$$\sqrt{3.01} - \sqrt{3} = \frac{3.01 - 3}{\sqrt{3.01} + \sqrt{3}} = \frac{0.01}{1.7349 + 1.7321} = 2.8844 \times 10^{-3}$$

也能达到结果具有5位有效数字的要求, 而这时  $\sqrt{3.01}$  和  $\sqrt{3}$  所需的有效位数只要5位, 远比直接相减所需有效位数(8位)要少。

同样由式(1.10)和式(1.11)两式可得乘、除法运算的误差公式

$$\varepsilon(x_1 x_2) \approx x_2^* \varepsilon(x_1) + x_1^* \varepsilon(x_2), \quad \varepsilon_r(x_1 x_2) \approx \varepsilon_r(x_1) + \varepsilon_r(x_2) \quad (1.13)$$

$$\varepsilon\left(\frac{x_1}{x_2}\right) \approx \frac{1}{x_2^*} \varepsilon(x_1) - \frac{x_1^*}{(x_2^*)^2} \varepsilon(x_2), \quad \varepsilon_r\left(\frac{x_1}{x_2}\right) \approx \varepsilon_r(x_1) - \varepsilon_r(x_2) \quad (1.14)$$

由式(1.13)的第一式可知, 当乘数很大时, 乘积的绝对误差可能很大, 应设法避免。由第二式可知, 近似值之积的相对误差等于相乘各因子的相对误差的代数和。由式(1.14)第一式可知, 当除数  $x_2^*$  的绝对值很小, 接近于零时, 商的绝对误差  $|\varepsilon\left(\frac{x_1}{x_2}\right)|$  可能会很大, 甚至造成计算机的“溢出”错误, 故应设法避免让绝对值太小的数作为除数。

综上分析可知, 大小相近的同号数相减, 乘数的绝对值很大, 以及除数接近于零等, 在数值计算中都应设法避免。

**例 1.4** 设已测得某长方形场地的长和宽的范围分别为  $L = 110 \pm 0.2 \text{ m}$ ,  $D = 80 \pm 0.1 \text{ m}$ , 求该场地的面积  $S$ , 并估算其绝对误差限和相对误差限。

**解** 由  $S = LD$  可求出面积  $S$  的近似值

$$S^* = 110 \times 80 = 8800 \text{ m}^2$$

由式(1.10)可计算  $S^*$  的绝对误差限, 由于

$$\varepsilon(S) \approx \left( \frac{\partial S}{\partial L} \right)^* \varepsilon(L) + \left( \frac{\partial S}{\partial D} \right)^* \varepsilon(D) = D^* \varepsilon(L) + L^* \varepsilon(D)$$

于是

$$|\varepsilon(S)| \leq |D^*| |\varepsilon(L)| + |L^*| |\varepsilon(D)| \leq 80 \times 0.2 + 110 \times 0.1 = 27 \text{ m}^2$$

因此,  $S^*$  的相对误差限为

$$\left| \frac{\varepsilon(S)}{S^*} \right| = \left| \frac{\varepsilon(S)}{8800} \right| \leq \frac{27}{8800} = 0.31\%$$

例 1.5 经过四舍五入得出  $x_1 = 6.1025, x_2 = 80.115$ , 求  $x_1 + x_2, x_1 x_2$  的绝对误差限。

解 由于  $|\varepsilon(x_1)| \leq \frac{1}{2} \times 10^{-4}, |\varepsilon(x_2)| \leq \frac{1}{2} \times 10^{-3}$ , 所以

$$\begin{aligned} |\varepsilon(x_1 + x_2)| &\approx |\varepsilon(x_1) + \varepsilon(x_2)| \leq |\varepsilon(x_1)| + |\varepsilon(x_2)| \\ &\leq \frac{1}{2} \times 10^{-4} + \frac{1}{2} \times 10^{-3} = 0.00055 \end{aligned}$$

## 1.4 算法的数值稳定性

通过前面对误差传播规律的分析, 同一问题当选用不同的算法时, 它们所得到的结果有时会相差很大, 这是因为运算中的舍入误差在运算过程中的传播常随算法而异。凡一种算法的计算结果受舍入误差的影响小者称它为数值稳定的算法。下面举几个例子来说明。

例 1.6 解方程

$$x^2 - (10^9 + 1)x + 10^9 = 0 \quad (1.15)$$

解 由韦达定理可知, 此方程的精确解为

$$x_1 = 10^9, \quad x_2 = 1$$

如果利用求根公式

$$x_{1,2} = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} \quad (1.16)$$

来编制计算机程序, 在字长为 8、基底为 10 的计算机上进行运算, 则由于计算机实际上采用的是规格化浮点数的运算, 这时

$$-b = 10^9 + 1 = 0.1 \times 10^{10} + 0.000\,000\,000\,1 \times 10^{10}$$

的第二项中最后两位数“01”, 由于计算机字长的限制, 在机器上表示不出来, 故在计算机对阶舍入运算时为

$$-b = 0.1 \times 10^{10} + 0.000\,000\,000\,1 \times 10^{10} = 0.1 \times 10^{10} = 10^9$$

$$\sqrt{b^2 - 4ac} = \sqrt{[-(10^9 + 1)]^2 - 4 \times 10^9} = 10^9$$

于是

$$x_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a} = \frac{10^9 + 10^9}{2} = 10^9$$

$$x_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a} = \frac{10^9 - 10^9}{2} = 0$$

这样算出的根  $x_2 = 0$  显然是严重失真的(因为精确解  $x_2 = 1$ )，这说明直接利用式(1.16)求解方程(1.15)是不稳定的。其原因是在于当计算机进行加、减运算时要对阶舍入计算，实际上受到机器字长的限制，绝对值相对小的数被大数“淹没”后就无法发挥其应有的影响，由此带来误差，造成计算结果的严重失真。这时，如要提高计算的数值稳定性，必须改进算法。在本例中算出的根  $x_1 = 10^9$  是可靠的，故可利用根与系数的关系式  $x_1 x_2 = \frac{c}{a}$  来计算  $x_2$ ，有

$$x_2 = \frac{c}{a} \cdot \frac{1}{x_1} = \frac{10^9}{1 \times 10^9} = 1$$

所得结果很好。这说明第二种算法有较好的数值稳定性。

**注：**在利用根与系数关系式求第二根时，必须先算出绝对值较大的一个根，然后再求另一个根，才能得到精度较高的结果。

### 例 1.7 试计算积分

$$I_n = \int_0^1 x^n e^{x-1} dx \quad (n = 1, 2, \dots)$$

**解** 由分部积分法可得

$$I_n = x^n e^{x-1} \Big|_0^1 - n \int_0^1 x^{n-1} e^{x-1} dx$$

因此，有递推公式  $I_n = 1 - nI_{n-1}$  ( $n = 1, 2, \dots$ )，其中  $I_1 = 1/e$ 。

用上面的递推公式，在字长为6，基底为10的计算机上，从  $I_1$  出发计算前几个积分值，其结果如表 1.1。

表 1.1

$k$	$I_k$
1	0.367 879
2	0.264 242
3	0.207 274
4	0.170 904
5	0.145 480
6	0.127 120
7	0.110 160
8	0.118 720
9	-0.068 480

被积函数  $x^n e^{x-1}$  在积分限  $(0, 1)$  区间内都是正值，积分值  $I_9$  取三位有效数字时的精确结果为 0.091 6，但上表中的  $I_9 = -0.068 480$  却是负值，与 0.091 6 相差很大。怎么会出现在这种现象？可分析如下。

由于在计算  $I_9$  时有舍入误差约为  $\varepsilon = 4.412 \times 10^{-7}$ ，且考虑以后的计算都不再另有舍入误

差。此  $\varepsilon$  对后面各项计算的影响为

$$\begin{aligned} I_2 &= 1 - 2(I_1 + \varepsilon) = 1 - 2I_1 - 2\varepsilon = 1 - 2I_1 - 2! \varepsilon \\ I_3 &= 1 - 3(I_2 + \varepsilon) = 1 - 3(1 - 2I_1) + 3! \varepsilon \\ I_4 &= 1 - 4[1 - 3(1 - 2I_1)] - 4! \varepsilon \\ &\vdots \\ I_9 &= 1 - 9[1 - 8(\cdots)] + 9! \varepsilon \end{aligned}$$

这样,算到  $I_9$  时产生的误差为  $9! \varepsilon \approx 0.6101$  就是一个不小的数值了。

可以改进算法来提高此例的数值稳定性,即将递推公式改写为

$$I_{n+1} = \frac{1 - I_n}{n}$$

从后向前进递推计算时,  $I_n$  的误差下降为原来的  $\frac{1}{n}$ ,因此只要  $n$  取得足够大,误差逐次下降,其影响就会越来越小。

由

$$I_n = \int_0^1 x^n e^{x-1} dx < \int_0^1 x^3 dx = \frac{1}{n+1}$$

可知:当  $n \rightarrow \infty$  时  $I_n \rightarrow 0$ 。因此可取  $I_{20} = 0$  作为初始值进行递推计算。

由于  $I_{20} \approx \frac{1}{20}$ ,故  $I_{20} = 0$  的误差约为  $\frac{1}{21}$ 。在计算  $I_{19}$  时误差下降到

$$\frac{1}{21} \times \frac{1}{20} \approx 0.0024$$

到计算  $I_{15}$  时,误差已下降到  $10^{-8}$  以下,结果如表 1.2。

这样得到的  $I_9 = 0.0916123$  已很精确了。可见经过改进后的新算法具有很好的稳定性。

表 1.2

$k$	$I_k$
20	0.000 000 0
19	0.050 000 0
18	0.050 000 0
17	0.052 777 8
16	0.055 719 0
15	0.059 017 6
14	0.062 732 2
13	0.066 947 7
12	0.071 773 3
11	0.077 352 3
10	0.083 877 1
9	0.091 612 3

**例 1.8** 对于小的  $x$  值, 计算  $e^x - 1$ 。

**解** 如果用  $e^x - 1$  直接进行计算, 其稳定性是很差的, 因为两个相近数相减会严重丢失有效数字, 产生很大的误差。因此得采用合适的算法来保证计算的数值稳定性。可以把  $e^x$  在点  $x = 0$  附近展开成幂级数:

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots$$

则可得

$$e^x - 1 = x \left( 1 + \frac{x}{2!} + \frac{x^2}{3!} + \dots \right)$$

按上式计算就有很好的数值稳定性。

通过以上这些例子, 可以知道算法的数值稳定性对于数值计算的重要性了。如无足够的稳定性, 将会导致计算的最终失败。为了防止误差传播、积累带来的危害, 提高计算的稳定性, 将前面分析所得的各种结果归纳起来, 得到数值计算中应注意如下几点:

- ① 应选用数值稳定的计算方法, 避开不稳定的算式。
- ② 注意简化计算步骤及公式, 减少误差的积累; 设法减少乘除法运算, 节约计算机的机时。
- 例如前面讲到过的用秦九韶算法计算多项式, 就是一个改变计算公式以减少运算次数的极好例子。
- ③ 应合理安排运算顺序, 防止参与运算的数在数量级相差悬殊时, 大数“淹没”小数的现象发生。
- ④ 应避免两相近数相减, 可用变换公式的方法来解决。
- ⑤ 绝对值太小的数不宜作为除数, 否则产生的误差过大, 甚至会在计算机中造成“溢出”错误。

## 习题 1

1. 下列各数都是对真值进行四舍五入后得到的近似值, 试分别写出它们的绝对误差限、相对误差限和有效数字的位数:

(1) $x_1^* = 0.024$	(2) $x_2^* = 100$
(3) $x_3^* = 57.50$	(4) $x_4^* = 8 \times 10^5$

2. 为了使  $\sqrt{11}$  的近似值的相对误差  $\leq 0.1\%$ , 问至少应取几位有效数字?

3. 求  $x$ , 使 20.345 和 20.346 作为它的近似值都具有 5 位有效数字。

4. 如用级数  $e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!}$  来求  $e^{-5}$  的值, 为使相对误差  $< 10^{-3}$ , 问至少需取几项?

5. 设  $x$  的相对误差限为 2%, 求  $x^n$  的相对误差限。

6. 为了使积分  $I = \int_0^1 e^{-x^2} dx$  的近似值  $I^*$  的相对误差不超过 0.1%, 问至少取几位有效数字?

7. 正方形的边长约为 100 cm, 问测量时误差最多只能到多少, 才能保证面积的误差不超过 1 cm<sup>2</sup>?

8. 已知  $y = P(x) = x^2 + x - 1150$ ,  $x = \frac{100}{3}$ ,  $x^* = 33$ , 计算  $y = P\left(\frac{100}{3}\right)$  及  $y^* = P(33)$ , 并