

SHUZHI SHIYAN YU
JISUANJI MONI

数值实验与 计算机模拟

李茜 张石峰 段海明◎著

$$\begin{aligned}S &\leftarrow \sin \theta = \sqrt{1 - \cos^2 \theta} \\D &\leftarrow R(K + 7) \\X_1 &\leftarrow X_0 + D \sin \theta \cos \phi \\Y_1 &\leftarrow Y_0 + D \sin \theta \sin \phi \\Z_1 &\leftarrow Z_0 + D \cos \theta\end{aligned}$$

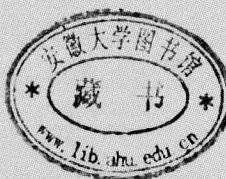


上海交通大学出版社
SHANGHAI JIAO TONG UNIVERSITY PRESS

SHUZHI SHIYAN YU
JISUANJI MONI

数值实验与 计算机模拟

李 茜 张石峰 段海明◎著



内 容 提 要

本书是作者多年来从事有关科学与工程项目研究及有关课程教学研究的成果。其中的一些研究专题取自作者承担的国家自然科学基金项目和工程研究项目,一些研究专题包含指导学士、硕士论文的内容。书中的数值试验部分是针对计算机工程系、数学系、物理系学生《数值分析》、《数值方法与程序设计》、《计算物理学》等课程的数值实验配套内容。

主要内容有:①数值实验的意义、有关算法及误差理论;②数值实验 10 个(含算法模型、实验内容及参考程序);③模拟研究专题 8 个(含模型概述、算法实现及参考程序)。全部参考程序均在 V-For6.5 上调试通过。

本书可供科技工作者做数值模拟和科学与工程计算时参考,也可供理工类专业教师从事教学与研究工作时借鉴。

图书在版编目(CIP)数据

数值实验与计算机模拟/李茜,张石峰,段海明著. —上海:
上海交通大学出版社,2013

ISBN 978 - 7 - 313 - 10383 - 3

I. ①数… II. ①李… ②张… ③段… III. ①数值计算-
实验-研究②数值计算-计算机模拟 IV. ①024

中国版本图书馆 CIP 数据核字(2013)第 233529 号

数值实验与计算机模拟

李 茜 张石峰 段海明 著

上海交通大学出版社出版发行

上海市番禺路 951 号 邮政编码 200030

电话: 64071208 出版人: 韩建民

凤凰数码印务有限公司印刷 全国新华书店经销

开本: 787mm×1092mm 1/16 印张: 10.75 字数: 254 千字

2013 年 10 月第 1 版 2013 年 10 月第 1 次印刷

ISBN 978 - 7 - 313 - 10383 - 3/O 定价: 25.00 元

版权所有 侵权必究

告读者: 如发现本书有印装质量问题请与印刷厂质量科联系
联系电话: 025 - 83657309

前　　言

计算机模拟也叫数值模拟。它以电子计算机为工具,通过数值计算和图像显示的方法,达到对科学与工程问题乃至自然界各类问题进行模拟、重现、预测的目的。数值计算方法和数值实验是它的理论和实践基础,高级程序设计语言是它的设计基础。本书选用科学与工程计算界常用的 Fortran 语言作为程序设计语言。

本书是作者多年来从事有关科学与工程项目研究及有关课程教学研究的成果。其中的一些研究专题取自作者承担的国家自然科学基金项目和工程研究项目,一些研究专题包含指导学士、硕士论文的内容。书中的数值试验部分是针对计算机工程系、数学系、物理系学生《数值分析》、《数值方法与程序设计》、《计算物理学》等课程的数值实验配套内容。在课程学习和论文撰写的过程中受到学生及青年教师的欢迎。本书主要内容有:①数值实验的意义、有关算法及误差理论;②数值实验 10 个(含算法模型、实验内容及参考程序);③模拟研究专题 8 个(含模型概述、算法实现及参考程序)。全部参考程序均在 V-For6.5 上调试通过。

肇庆工商职业技术学院的李茜老师负责编写了本书第 1 篇的 1.1~1.3,第 2 篇的实验 2.1~2.4,第 3 篇的专题 3.3、3.4;张石峰老师负责编写了第 1 篇的 1.4、1.5,第 2 篇的实验 2.6~2.8,第 3 篇的专题 3.5~3.8;新疆大学段海明老师负责编写了第 2 篇的实验 2.9、2.10,第 3 篇的专题 3.1、3.2;新疆师范大学王林香老师负责编写了第 2 篇的实验 2.5,第 3 篇的专题 3.6。本书由李茜、张石峰负责统稿。由于编者水平所限,难免不足之处,恳请读者提出宝贵意见和建议。

本书可供科技工作者做数值模拟和科学与工程计算时的参考,也可供理工类专业教师从事教学与研究工作时借鉴。

本书的出版得到上海交通大学出版社和肇庆工商职业技术学院科研处的大力支持,在此表示衷心感谢!

编　者

2013 年 6 月于广东肇庆 北岭山麓

目 录

第 1 篇 引论	1
1. 1 数值实验及其步骤	1
1. 2 数值实验的重要性	1
1. 3 误差的来源与误差的基本概念	3
1. 3. 1 误差的来源与类型	3
1. 3. 2 绝对误差与绝对误差限	4
1. 3. 3 相对误差与相对误差限	5
1. 3. 4 有效数字	5
1. 4 数值计算中需要注意的问题	7
1. 4. 1 避免两个相近的数相减	7
1. 4. 2 防止大数“吃掉”小数	8
1. 4. 3 注意简化计算步骤,减少运算次数	9
1. 5 程序设计中应该注意的问题	10
1. 5. 1 对程序设计的一些建议	10
1. 5. 2 程序设计中的一些经验	11
 第 2 篇 数值实验	17
2. 1 Fortran 语言程序的调试运行过程	17
2. 1. 1 Fortran 语言简介	17
2. 1. 2 Fortran 程序在 V-F 环境下的调试运行过程	17
2. 1. 3 调试运行例程序	18
2. 1. 4 实验与思考	19
2. 2 实验数据处理 I 与插值法	20
2. 2. 1 实验目的	20
2. 2. 2 算法描述	20
2. 2. 3 实验内容	20
2. 2. 4 参考程序	20
2. 3 实验数据处理 II 与曲线拟合	25

2.3.1 实验目的	25
2.3.2 算法描述	25
2.3.3 实验内容	26
2.3.4 参考程序	26
2.4 解线性方程组的迭代法	30
2.4.1 实验目的	30
2.4.2 算法描述	30
2.4.3 实验内容	31
2.4.4 参考程序	31
2.5 解线性方程组的消去法	33
2.5.1 实验目的	33
2.5.2 算法描述	34
2.5.3 实验内容	34
2.5.4 参考程序	34
2.6 非线性方程求根	38
2.6.1 实验目的	38
2.6.2 算法步骤	38
2.6.3 实验内容与要求	38
2.6.4 参考程序	39
2.7 矩阵特征值的计算	43
2.7.1 实验目的	43
2.7.2 幂法算法描述	43
2.7.3 实验内容	44
2.7.4 参考程序	44
2.8 常微分方程(组)的数值解法	49
2.8.1 实验目的	49
2.8.2 算法描述	50
2.8.3 实验内容	50
2.8.4 参考程序	50
2.8.5 思考练习	56
2.9 随机数的产生与检验	56
2.9.1 实验目的	56
2.9.2 算法描述	56
2.9.3 实验内容	58
2.9.4 参考程序	58

2.9.5 思考练习	59
2.10 蒙特卡罗方法的应用	60
2.10.1 实验目的	60
2.10.2 算法描述	60
2.10.3 实验内容	64
2.10.4 参考程序	64
 第3篇 研究专题	70
3.1 无规网格点方法求解 LJ 团簇的基态几何结构*	70
3.1.1 Lennard-Jones (LJ) 团簇简介	70
3.1.2 LJ 团簇的势函数描述	70
3.1.3 无規格点方法、模型构造与优化	71
3.1.4 最速下降法	72
3.1.5 无規格点方法求解 LJ 团簇的基态几何结构	72
3.1.6 程序使用说明	72
3.1.7 源程序	73
3.2 蒙特卡罗模拟退火方法求解 LJ 团簇的基态几何结构*	81
3.2.1 蒙特卡罗方法	81
3.2.2 Metropolis 蒙特卡罗模拟退火方法	81
3.2.3 蒙特卡罗模拟退火方法求解 LJ 团簇的基态几何结构	81
3.2.4 程序使用说明	82
3.2.5 源程序	83
3.3 中子星星体结构的计算	88
3.3.1 星体结构微分方程(组)	88
3.3.2 广义相对论效应	89
3.3.3 物态方程及微分方程(组)的化简	89
3.3.4 中子星结构特性计算	90
3.3.5 参考程序	90
3.4 逾渗模型的研究	98
3.4.1 逾渗理论的描述	99
3.4.2 二维正方逾渗模型的临界点研究	101
3.4.3 三维正方逾渗模型的临界点研究	108
3.5 有限扩散凝聚模型(DLA)	116
3.5.1 DLA 模型简介	117
3.5.2 DLA 生长过程的 M-C 模拟	118

3.5.3 参考程序	122
3.6 低能离子注入植物种子的深度-浓度分布的模拟*	129
3.6.1 有关离子注入的 LSS 理论	129
3.6.2 近似假设	130
3.6.3 参考程序	131
3.7 中子输运问题的 M-C 模拟	136
3.7.1 中子输送问题	136
3.7.2 中子穿透平板模型	136
3.7.3 直接法与加权法模拟	137
3.7.4 参考程序	139
3.8 地下水渗流问题中最紧凑存储的有限单元方法*	145
3.8.1 地下水渗流运动的有限元方程	145
3.8.2 按结点集成有限元方程组及一维最紧凑存储的方法	147
3.8.3 非线性问题的处理	150
3.8.4 参考程序	150
附录 1 FORTRAN 语句表	154
附录 2 FORTRAN 函数表	156
参考文献	159

注:以上加 * 者为国家自然科学基金项目内容的选题

第1篇 引论

“数值试验”是《数值分析》、《数值方法与程序设计》、《高级语言程序设计》、《计算物理学》等课程不可缺少的实践环节。学生通过“数值试验”，逐步熟悉并掌握数值分析与程序设计的基本理论与步骤。

1.1 数值实验及其步骤

数值实验不同于传统的物理实验(实验室实验)。广义地说,数值实验是用计算机及物理数学的方法来模拟、了解无法解析求出的客观事物的状态和过程的实践;狭义地说,数值实验就是在计算物理学中与各种理论方法相对应的程序实现过程及结果分析。数值实验往往要从数学模型得到离散化的数值模型,所得结果也往往是离散的数值结果,这是它属于又有别于数学实验的地方。

数值实验的一般步骤是:物理模型→数学模型→数值模型→程序设计→模型校正→模型运行→数值(或图形)结果→结果分析。

与传统实验比较,前3步建立模型,相当于确定实验对象(有时也可以是两步:直接从物理模型得到数值模型,例如天然差分方法);程序设计相当于仪器选配;模型校正相当于仪器校准;模型运行相当于测量过程。

要特别重视模型校正与结果分析。没有一个经反复调试验证的模拟模型,不可能得到符合客观的实验结果,这需要反复对前面的程序设计、数值模型以至数学模型、物理模型进行校验,并对模型中的参数进行辨识,参数辨识可以用最优化方法进行逼近。

结果分析也是数值实验的重要环节,要充分利用数值实验易于重复的特点对数值结果进行多方位的试验分析研究。

1.2 数值实验的重要性

数值实验是计算物理学的重要研究手段。数值实验可以解决或补充解决实验物理和理论物理难于解决甚至无法解决的问题。例如在天体力学中的多体问题、湍流理论的发展、金属裂缝的传播、油藏多相流模拟等问题中,只有少部分能够用解析方法解决,大部分需要通过物理实验和数值实验的方法来解决;核爆炸、核电站试验中,数值实验是更经济有效的实

验手段;通过数值实验设计飞机外型的结果,机翼的阻力比用机械实验设计的结果要小40%。

物理模型的正确性及误差,数学模型描述物理模型的正确性及误差,数值算法求解数学模型的正确性及误差,计算机程序实现数值算法的正确性及误差……都可以通过数值实验的结果来分析比较、判定和检验。因此,数值实验是学习掌握计算物理学的重要基本功。

在数值实验中,数值方法的选取、离散区间的选取、计算的方向和顺序常常是重要的甚至是决定性的。如以下实例:

[例 1.1] 由常用函数幂级数展开式可得

$$\ln 2 = 1 - \frac{1}{2} + \frac{1}{3} - \frac{1}{4} + \cdots + (-1)^n \frac{1}{n} + \cdots$$

$$\text{若记 } x_n = \sum_{k=1}^n \frac{(-1)^{k-1}}{k}$$

则 x_n 是逼近 $\ln 2$ 的数列,并且有估计式

$$|x_n - \ln 2| < \frac{1}{n+1}$$

若需 $|x_n - \ln 2| < \epsilon$, 要取 $\frac{1}{n+1} < \epsilon$ 。当 $\epsilon = \frac{1}{2} \times 10^{-4}$ 时, $n \approx 2000$; 当 $\epsilon = \frac{1}{2} \times 10^{-5}$ 时, $n \approx$

20000! 即要将 $\ln 2$ 的有效数字从 4 位提高到 5 位,则求和的项数要从 2000 项提高到 20000 项!

实际上,由于计算机舍入误差积累的影响,一共要计算 $n = 99274$ 次! 这表明 x_n 是一个收敛很慢的数列。由于舍入误差的积累随计算次数的膨胀,实际上得不到满足 $\epsilon = \frac{1}{2} \times 10^{-8}$

条件的 n 。

但若将 x_n 做简单的组合变换

$$x'_n = x_n - (x_n - x_{n-1})^2 / (x_n - 2x_{n-1} + 2x_{n-2})^2, \quad n=3,4,\dots \quad (2)$$

x'_n 为新的逼近 $\ln 2$ 的数列, 数值实验表明 x'_n 以很快的速度收敛到 $\ln 2$ 。精度达到 $1/2 \times 10^{-5}$ 时, 只需要计算 30 项! 精度达到 $1/2 \times 10^{-6}$ 时, 只需要计算 64 项! 比原来直接计算 x_n 减少了千万倍!

看来理论上可行的算法,也应经过数值实验来检验其优劣。

[例 1.2] 在计算机辅助设计(CAD)中常用折线去逼近曲线,以 $y = \sqrt{x}, 0 \leq x \leq 1$ 为例。在 $[0,1]$ 区间上取 $0 = x_0 \leq x_1 \leq \cdots \leq x_n = 1$, 在每个小区间 $[x_j, x_{j+1}]$ 对函数 \sqrt{x} 建立线性插值函数 S_j , 则 $y = S_j(x), x_j \leq x \leq x_{j+1}, j=0,1,\dots,n-1$, 即是逼近曲线 $y = \sqrt{x}$ 的折线段。在 $x_j \leq x \leq x_{j+1}$ 区间上的最大误差为

$$\max |x^{1/2} - S_j(x)| = \frac{1}{4} (x_{j+1}^{1/2} - x_j^{1/2}) / (x_{j+1}^{1/2} + x_j^{1/2})$$

问题是:当 n 固定,为提高精度,要求 $E = \max |x^{1/2} - S_j(x)|$ 尽可能小,问 x_j 应如何选取?

做法1:对 $[0,1]$ 区间做等距划分,即取 $x_j = \frac{j}{n}$, $j=0,1,2,\dots,n$ 。则

$$E = \frac{1}{4\sqrt{n}}$$

做法2:对 $[0,1]$ 区间做不等距划分,即取 $x_j = \left(\frac{j}{n}\right)^4$, $j=0,1,2,\dots,n$ 。则

$$E = \frac{1}{4n^2} \cdot \frac{4j^2 + 4j + 1}{2j^2 + 2j + 1} \leqslant \frac{1}{2n^2}$$

若要求 $E < 10^{-4}$,则在做法1中 n 要取 25×10^6 ,即用2500万段折线才能达到万分之一的精度,计算量非常浩大!若采用做法2,则只要取 $n=71$,即可达到同样的精度。

仅仅选取不同的插值点,也能成千上万倍地减少计算量!

例1.3 考虑如下积分

$$I_n = e^{-1} \int_0^1 x^{-n} e^x dx \quad \text{取 } n = 0, 1, 2, \dots, 7, \text{ 对应 } 8 \text{ 个积分。}$$

它们满足递推关系

$$I_n = 1 - n I_{n-1}$$

若算出 I_0 后用递推公式去求 I_1, I_2, \dots, I_7 ,则计算误差递增,计算结果不可靠。若用反方向的递推关系 $I_{n-1} = (1 - I_n)/n$,则计算误差逐步削弱,可得到可靠的结果。

看来计算的方向和顺序也会影响计算精度、收敛速度甚至计算的稳定性。

1.3 误差的来源与误差的基本概念

1.3.1 误差的来源与类型

一个物理量的真实值与我们计算出的数值往往不相等,我们称其为误差。引起误差的原因是多方面的,为了具体说明误差的来源,我们先来分析一下解决实际问题的过程:



从以上过程我们可以看出,实际问题与计算结果存在着以下几种误差:

1) 模型误差

数学模型是从实际问题经抽象和简化,并忽略一些次要因素得到的。例如用 $s = \frac{1}{2}gt^2$

(其中 g 是重力加速度)描述地球上某一质点自由落体运动规律时会出现误差(因为忽略了空气阻力等一些次要因素),我们把这种数学模型与实际问题之间出现的误差称为模型误差。

若用 $s = f(t)$ 表示自由落体的真实运动规律,那么 $f(t) - \frac{1}{2}gt^2$ 即为质点自由落体运动

数学模型的模型误差。

2) 参数误差

在给出的数学模型中往往涉及一些根据观测得到的物理量,如电压、电流、温度、长度等,而观测难免不带来误差,观测值与真实值之间的误差称为参数误差或观测误差。

例如,取重力加速度为 $g^* \approx 9.8 \text{m/s}^2$,则 $g - g^*$ 就是参数误差。

3) 截断误差

在计算中常常遇到只有通过无限计算过程才能得到结果的情况,但实际计算时,只能用有限过程来计算。这种用有限过程代替无限过程的误差称为截断误差,而这种误差是由计算方法本身引起的,因此也称为方法误差。

例如,考虑用 Taylor 级数求定积分 $\int_0^1 \frac{\sin x}{x} dx$ 的近似值。由 $\sin x$ 的 Taylor 展开式得

$$\frac{\sin x}{x} = \frac{1}{x} \left(x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots \right)$$

若从第二项后“截断”,则有

$$\int_0^1 \frac{\sin x}{x} dx \approx \int_0^1 \frac{1}{x} \left(x - \frac{x^3}{3!} \right) dx = \int_0^1 \left(1 - \frac{x^2}{6} \right) dx = 1 - \frac{1}{18} = \frac{17}{18}$$

这时 $\int_0^1 \frac{\sin x}{x} dx - \frac{17}{18}$ 即为截断误差。

4) 舍入误差

在计算中遇到的数据可能位数很多,也可能是无穷小数,如 $\sqrt{2}$ 、 $1/3$ 等,但计算时只能对有限位数进行计算,因此往往进行四舍五入,这样产生的误差称为舍入误差。

少量的舍入误差是微不足道的,但在电子计算机上完成了千百万次运算之后,舍入误差的累积有时可能是十分惊人的。

在上述讨论的误差来源中,前两种误差是客观存在的,后两种误差是由数值计算方法所引起的。

我们的数值实验主要考虑从数学模型去计算离散的数值结果,因此主要涉及后两种误差。

1.3.2 绝对误差与绝对误差限

定义 1 设 x 是准确值, x^* 是它的一个近似值, 则 $x - x^*$ 为 x^* 的绝对误差,简称误差,记作 e^* ,即

$$e^* = x - x^* \quad ①$$

例如,用 1.414 近似 $\sqrt{2}$,其绝对误差为:

$$\sqrt{2} - 1.414 = 1.414213\dots - 1.414 = 0.000213\dots$$

通常我们无法知道准确值 x ,因而也不可能知道误差 e^* 的准确值。但我们很容易得到 e^* 的取值范围。例如用 1.414 作为 $\sqrt{2}$ 的近似值,其绝对误差不会超过 0.0003,因此我们可以给出绝对误差的上限。

定义2 设 x 为准确值, x^* 是它的一个近似值,称 x^* 的绝对误差的绝对值的上限 ϵ^* 为 x^* 的绝对误差限,简称误差限,即

$$|\epsilon^*| = |x - x^*| \leq \epsilon^* \quad (2)$$

显然,如果 ϵ^* 是 x 的近似值 x^* 的绝对误差限,那么 x 仅位于区间 $[x^* - \epsilon^*, x^* + \epsilon^*]$ 之间,在工程技术上常用 $x = x^* \pm \epsilon^*$ 表示。例如用毫米刻度的直尺测量某一长度为 x 的物体,测得其长度的近似值为 $x^* = 52\text{mm}$,由于直尺以毫米为刻度,其误差不超过 0.5mm ,故 $x = 52 \pm 0.5\text{mm}$ 。

1.3.3 相对误差与相对误差限

在许多情形下,绝对误差限并不能完全刻划一个数的近似精确程度。例如,比较 $x = 10 \pm 1$ 和 $y = 10000 \pm 10$ 两种情形。从绝对误差限来看, y^* 的绝对误差限是 x^* 的绝对误差限的十倍;但从实际情况来看, y^* 的精确程度要高于 x^* 的精确程度。因此,一个近似值的精确程度不仅与绝对误差限有关,而且还与其本身的大小有关。由此我们给出以下定义:

定义3 设 x 为准确值, x^* 是它的一个近似值;称比值 $\frac{\epsilon^*}{x^*}$ 为近似值 x^* 的相对误差,记作 ϵ_r^* ,即

$$\epsilon_r^* = \frac{\epsilon^*}{x^*} = \frac{x - x^*}{x^*} \quad (3)$$

与绝对误差限的概念类似,我们引入相对误差限的概念:

定义4 设 x 为准确值, x^* 是它的一个近似值,称 x^* 的相对误差 ϵ_r^* 的绝对值的上界 ϵ_r^* 为 x^* 的相对误差限,即

$$|\epsilon_r^*| \leq \epsilon_r^* \quad (4)$$

相对误差和相对误差限都是无量纲的,常用百分数表示。例如, $x = 10 \pm 1$ 和 $y = 10000 \pm 10$ 的近似值 $x^* = 10$ 和 $y^* = 10000$ 的相对误差限分别为

$$\epsilon_r^*(x) = \frac{1}{10} = 10\%, \quad \epsilon_r^*(y) = \frac{10}{10000} = 0.1\%$$

x^* 的相对误差限是 y^* 的相对误差限的 100 倍,在这种意义上说, y^* 的精度比 x^* 的精度高得多,这是符合实际情况的。

由定义2与定义4,可以得到绝对误差限与相对误差限之间的关系为

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

1.3.4 有效数字

在实际计算中,经常按四舍五入原则取近似数。例如:

$$\sqrt{200} = 14.142\dots \approx 14.1, \quad \epsilon_1^* < 0.05$$

$$\lg 2 = 0.30102\dots \approx 0.301, \quad \epsilon_2^* < 0.0001$$

$$e^{-5} = 0.0067379\dots \approx 0.00674, \quad \epsilon_3^* < 0.000003$$

它们的误差都不会超过末位数字的一半,即

$$|\sqrt{200} - 14.1| < \frac{1}{2} \times 10^{-1} = \frac{1}{2} \times 10^{2-3}$$

$$|\lg 2 - 0.301| < \frac{1}{2} \times 10^{-3} = \frac{1}{2} \times 10^{0-3}$$

$$|e^{-5} - 0.00674| < \frac{1}{2} \times 10^{-5} = \frac{1}{2} \times 10^{-2-3}$$

定义 5 设 x 为准确值, x^* 是它的一个近似值, 若将 x^* 表示成

$$x^* = \pm 0.a_1a_2\cdots a_n \times 10^m \quad (5)$$

其中 m, n 为整数, a_1, a_2, \dots, a_n 为 $0 \sim 9$ 之间的数, 且 $a_1 \neq 0$, 并满足关系式

$$|x - x^*| < \frac{1}{2} \times 10^{m-n} \quad (6)$$

则称 x^* 具有 n 位有效数字。

[例 1.4] 写出 $1/27$ 的具有一位、两位、三位和四位有效数字的近似值。

解 由于 $1/27 = 0.037037037\dots$

则按照定义得到一位、两位、三位和四位有效数字的近似值分别为 $0.04, 0.037, 0.0370, 0.03704$ 。

注意: 0.037 与 0.0370 两者之间是有差别的, 前者具有两位有效数字, 其误差不超过 $\frac{1}{2} \times 10^{-3}$,

而后者具有三位有效数字, 其误差不超过 $\frac{1}{2} \times 10^{-4}$ 。

定义 5 给出了绝对误差限与有效数字之间的关系, 下面我们给出相对误差限与有效数字之间的关系。

定理 1 设形如(5)式的近似值 x^* 具有 n 位有效数字, 则其相对误差限为

$$|e_r^*| \leq \frac{1}{2a_1} \times 10^{-(n-1)} \quad (7)$$

其中 a_1 是 x^* 的第一位有效数字。

证明: 由(5)式知,

$$|x^*| = 0.a_1a_2\cdots a_n \times 10^m \geq 0.a_1 \times 10^m$$

又由(6)式, 有

$$|e_r^*| = \frac{|x - x^*|}{|x^*|} \leq \frac{\frac{1}{2} \times 10^{m-n}}{0.a_1 \times 10^m} = \frac{1}{2a_1} \times 10^{-(n-1)}$$

定理证毕。

由定理 1 可以看出, 有效数字越多, 相对误差限越小。

并且很容易得到例 1.4 的相对误差限分别为

$$\epsilon_r^* = \frac{1}{2 \times 4} \times 10^{-(1-1)} = 12.5\%$$

$$\epsilon_r^* = \frac{1}{2 \times 3} \times 10^{-(2-1)} = 1.7\%$$

$$\epsilon_r^* = \frac{1}{2 \times 3} \times 10^{-(3-1)} = 0.17\%$$

$$\epsilon_r^* = \frac{1}{2 \times 3} \times 10^{-(4-1)} = 0.017\%$$

[例 1.5] 为使 $\sqrt{200}$ 的近似值的相对误差不超过 0.1%，问要取几位有效数字？

解 由⑦式，只需求出满足 $\frac{1}{2a_1} \times 10^{-(n-1)} \leq 0.1\%$ 的 n 即可。显然，近似值的第一位有效数字为 $a_1=1$ 。因此，由

$$\frac{1}{2} \times 10^{-(n-1)} \leq 0.1\%$$

可得 $n \geq 4$ ，于是 $\sqrt{200} \approx 14.14$

1.4 数值计算中需要注意的问题

1.4.1 避免两个相近的数相减

在数值计算中，两个相近的数相减，有时会严重损失有效数字，因而导致很大的相对误差。

例如， $x=1.5846$, $y=1.5839$ 都具有 5 位有效数字，但 $x-y=0.0007$ 只有一位有效数字。

为了更清楚地了解相近数相减的危害，我们看例 1.6。

[例 1.6] 设 $x=18.496$, $y=18.493$, 取四位有效数字计算 $x-y$ 的近似值，并估计其相对误差。

解 取 $x^*=18.50$, $y^*=18.49$, 则

$$x^*-y^*=18.50-18.49=0.01$$

而

$$x-y=18.496-18.493=0.003$$

其相对误差为

$$|\epsilon_r^*| = \left| \frac{(x-y)-(x^*-y^*)}{x^*-y^*} \right| = \left| \frac{0.003-0.01}{0.01} \right| = 70\%$$

由此例可以看出相对误差变得很大。

现在从理论上分析一下两数之差的相对误差。设 x^* 和 y^* 分别为 x 和 y 的近似值，因而 x^*-y^* 是 $x-y$ 的近似值，其相对误差应满足

$$|\epsilon_r^*| = \frac{|(x-y)-(x^*-y^*)|}{|x^*-y^*|} \leq \frac{|x-x^*|}{|x^*-y^*|} + \frac{|y-y^*|}{|x^*-y^*|}$$

$$= \frac{|x^*|}{|x^* - y^*|} |e_r^*(x)| + \frac{|y^*|}{|x^* - y^*|} |e_r^*(y)|$$

当 x^* 和 y^* 非常接近时, $|x^* - y^*|$ 很小, 使得相对误差 $|e_r^*|$ 变得很大。

为了避免上述情况的发生, 可以利用恒等式改变计算方法, 减少有效数字的损失。例如, 当 x 和 y 是非常接近的正数时, 有

$$\lg x - \lg y = \lg(x/y)$$

当 x 很大时, 有

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

在无恒等式可用的情况下, 可选用 Taylor 展开式, 如当 $f(x) \approx f(x^*)$ 时, 有 $f(x) - f(x^*) = f'(x^*)(x - x^*) + 1/2f''(x^*)(x - x^*)^2 + \dots$

取右端的有限项近似计算。例如, 当 $x^* \approx 0$ 时, 计算 $e^x - 1$, 可用

$$e^x - 1 = x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} + \dots$$

如果无法改变算式, 则可采用增加变量有效位数的方法: 在计算机上采用双倍字长(定义变量为双精度类型)运算, 以保证必须的精度。

1.4.2 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, 又要作对阶处理, 如不注意运算次序与方法, 就可能出现大数“吃掉”小数的现象, 使小数不能发挥其作用, 影响计算结果的准确性。

[例 1.7] 求方程

$$x^2 + (\alpha + \beta)x + 10^9 = 0$$

的根, 其中 $\alpha = -10^9$, $\beta = -1$

解 显然, 方程的两个根为 $x_1 = 10^9$, $x_2 = 1$, 如果我们用 8 位数字计算机, 使用二次方程的求根公式

$$x_1, x_2 = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

进行计算, 则

$$\begin{aligned} -b &= -(\alpha + \beta) = 10^9 + 1 \\ &= 0.10000000 \times 10^{10} + 0.10000000 \times 10^1 \\ &= 0.10000000 \times 10^{10} + 0.000000001 \times 10^{10} \\ &\doteq 0.10000000 \times 10^{10} \quad (\text{第 } 9, 10 \text{ 位舍去}) \\ &= 10^9 = -\alpha \end{aligned}$$

(上式中的符号 \doteq 表示机器中的相等)那么有:

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

所以

$$x_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a} = \frac{10^9 + 10^9}{2 \times 1} = 10^9$$

$$x_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a} = \frac{10^9 - 10^9}{2 \times 1} = 0$$

实际上, x_2 应等于 1。可以看出 x_2 的误差太大。原因是在作加、减法运算过程时要“对阶”, 因而小数 1 在对阶过程中, 被大数 10^9 吃掉了。从上述计算可以看出 x_1 的计算是可靠的, 而 x_2 的计算是不可靠的。我们利用两根间的关系 $x_1 \cdot x_2 = c/a$, 求出

$$x_2 = \frac{c}{ax_1} = \frac{10^9}{1 \times 10^9} = 1$$

此方法是可靠的。

[例 1.8] 计算定积分 $\int_N^{N+1} \ln x dx$, 其中 N 是某个很大的正整数。

解 由分部积分公式可得

$$\int_N^{N+1} \ln x dx = (N+1) \ln(N+1) - N \ln N - 1 = -1$$

(小数 1 被大数 N 吃掉了)而实际上由 $y = \ln x$ 的单调性得到不等式

$$\int_N^{N+1} \ln x dx > \int_N^{N+1} \ln N dx = \ln N$$

可见计算结果严重失真。

为了提高计算精度, 把上式改为

$$\begin{aligned} \int_N^{N+1} \ln x dx &= (N+1) \ln(N+1) - N \ln N - 1 \\ &= N[\ln(N+1) - \ln N] + \ln(N+1) - 1 \\ &= N \ln \frac{N+1}{N} + \ln(N+1) - 1 \\ &= N \ln \left(1 + \frac{1}{N}\right) + \ln(N+1) - 1 \\ &= N \left(\frac{1}{N} - \frac{1}{2N^2} + \frac{1}{3N^3} - \frac{1}{4N^4} + \dots\right) + \ln(N+1) - 1 \\ &= \ln(N+1) - \frac{1}{2N} + \frac{1}{3N^2} - \frac{1}{4N^3} + \dots \end{aligned}$$

此时即使小数 1 被大数 N 吃掉, 产生的误差也是非常小的。

1.4.3 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, 不但可以节省计算机的计算时间, 还能减少舍入误差。这是数值计算中必须遵守的原则, 也是数值实验要研究的重要内容。

[例 1.9] 计算 x^{31} 的值。

解 若将 x 的值逐个相乘, 那么要做 30 次乘法, 若写成